Machine Learning Summer School Kyoto 2012 Graph Mining Chapter 1: Data Mining

National Institute of Advanced Industrial Science and Technology

Koji Tsuda

National Institute of Advanced Industrial Science and Technology (AIST)

- Research institute under METI for 6 fields: Life Science, Informatics, Environment and Energy, Nanotechnology and Materials, Geological Survey, Standard
- 2500 scientists, 700 administrative staff, 5200 scientists from outside
- Since 2001

# Computational Biology Research Center

#### Odaiba, Tokyo

- Developing Novel Methods and Tools for
  - Genome Informatics
  - Molecular Informatics
  - Cellular Informatics

 Diverse Collaboration with Companies and Universities



# Koji Tsuda: Short Bio

1998 PhD in Kyoto Univ, Join ETL 2000 GMD FIRST (Berlin, Germany) 2001 Join CBRC/AIST 2003-2004, 2006-2008 Max Planck Institute for Biological Cybernetics, Tuebingen, Germany 2009 Back to CBRC, Machine Learning **Group Leader** 

# About this lecture

 How to extract knowledge from structured data
 Itemset mining, tree mining, graph mining

"Reverse Search Principle"
 Learning from structured data

Kernels for structured data

# Chapter 1 (Data Mining)



# Agenda 2 (Learning from Structured Data)

1. Preliminaries
2. Graph Clustering by EM
3. Graph Boosting
4. Regularization Paths in Graph Classification
5. Itemset Boosting for predicting HIV drug resistance

# Agenda 3 (Kernel)

1. Kernel Method Revisited
2. Marginalized Kernels (Fisher Kernels)
3. Marginalized Graph Kernels
4. Weisfeiler-Lehman kernels
5. Reaction Graph kernels
6. Concluding Remark

# Part 1: Structural Data in Biology

# **Biological Sequences**

DNA sequences (A,C,G,T) Gene Finding, Splice Sites RNA sequences (A,C,G,U) MicroRNA discovery, etc. Amino acid sequences (20 symbols) Remote homolog detection, Fold recognition etc.

## Structures hidden in sequences (I)

#### Exon/intron of DNA (Gene)



# Structures hidden in sequences (II)



 It is crucial to infer hidden structures and exploit them for classification
 Biological Graphs



hydrogens

# Gene Expression Data

Measurement of many mRNAs in the cell Rough estimate of amount of proteins Time-series or not Snapshot of the underlying dynamic system

			•				ŕ.	1.	1		Y		٠	٠				1	•			F.					٠	1				
							ċ.			0	1	٠		٠		۰.																
			•						5	.64	X.		•	٠			۰.															
		٠	•							2			•	•						•			-									•
		٠				1	Ζ.	1		v			2																	٠		•
		٠	•				۰.																-				٠	٠			1	
. 4			•																											٠		
			•																							٠						
			•0									2																		•		
1		٠	•																								٠					
. 5	٠	٠																														
. +													٠													٠	٠	٠				
			•	1											1				٠										•	٠		
- 23		٠	•		e,	e,																				٠						
1		٠						٠																								
		٠	•																								٠					
													٠								٠											
٠			•																		٠										• •	
. *		٠		•							•			٠	٠				•			•										
						Ψ.															•											
•			•	•	•	5.1									٠															٠		
				•										•	٠						•				٠			٠				
1	۰,	2		•		η,							•								•											
			۰.	۰.										•							•										5	
				•	•									٠	٠				٠					• •								
																		۰.				۰.									• •	
•		٠	•		•		• •								٠															•		
			•	•								•		٠			•				•							٠				
											٠		٠								٠											
	2													٠		٠			٠										2			
			•			•				٠	٠	۰								۰.		٠							•			
						•																	5.					۰.				
			۰.				c.	•			•									٠									2			
				•						٠							٠	٠			٠						٠					
									.*		٠		٠		1						٠											
	2													٠		٠			٠								2		P	2		
										1		٠																	1			
																												1				
																												2				

# **Biological Networks**

Protein-protein physical interaction
 Metabolic networks
 Gene regulatory networks
 Network induced from sequence similarity

Thousands of nodes (genes/proteins)
 100000s of edges (interactions)

# **Physical Interaction Network**

# Metabolic Network



00020 3/19/04

17

# Many possible prediction problems..

#### Given Data

#### **Predicted Property**



structure (3D coordinates of the atoms)

function (e.g., according to GO or MIPS)

interactions (with other proteins, DNA or metabolites)

localization (e.g., compartment)

# Part 2: Itemset mining

## **Data Mining**

- A formal study of efficient methods for extracting interesting rules and patterns from massive data
- Frequent itemset mining (Agrawal and Srikant 1994)
- Closed pattern mining
- Structured data mining (Sequence, Trees, and Graphs)





#### **Definitions: Database**

- A set  $\Sigma = \{ 1, ..., n \}$  of items (elements)
- Transaction database

 $L = \{1, 2, 3, 4\}$ 

- A set  $T = \{ t_1, ..., t_m \}$  of subsets of  $\Sigma$
- Each subset  $t \subseteq \Sigma$  is called a transaction

#### Alphabet of items

id	transaction
t1	1, 3
t2	2, 4
t3	1, 2, 3, 4
t4	1, 2, 4

#### **Definitions: Frequent sets**

- Itemset X appears in transaction t:  $X \subseteq t$
- Occurrence of X in database T:  $Occ(X, T) = \{ t \in T : X \subseteq t \}$
- Frequency of X: Fr(X, T) = |Occ(X, T)|
- Minimum support (minsup):  $0 \leq \sigma \leq |\mathbf{T}|$
- X is frequent in T if  $Fr(X, T) \ge \sigma$ .

 $I = \{1, 2, 3, 4\}$ 

Alphabet of items Transaction database

#### **Occurrences and frequencies** of itemsets

 $Occ(3, T) = \{t1, t3\}$ Fr(3, T) = 2 $Occ(24, T) = \{t2, t3, t4\},\$ Fr(24, T) = 3•23

id	transaction
t1	1, 3
t2	2, 4
t3	1, 2, 3, 4
t4	1, 2, 4

## Market Basket Data

Popular application of itemset mining

Business and Market data analysis

Transaction Data	ID	Chips	Mustard	Sausage	Softdrink	Beer
of purchase	001	1	0	0	0	1
•	002	1	1	1	1	1
$\rightarrow$	003	1	0	1	0	0
• a transaction	004	0	0	1	0	1
or a "basket"	005	0	1	1	1	1
of a basket	006	1	1	1	0	1
	007	1	0	1	1	1
	008	1	1	1	0	0
	009	1	0	0	1	0

#### •Meaning of the transaction 003

"Custmer 003 bought Chips and Sausage together in his basket"







## **Association Rule Mining**

- Confidence: Supp(A U B)/ Supp(A)
   Probability of B, Given A
- What item is likely to be bought when A is bought
- Search: large support, confidence large
- Post-processing of itemset mining

## **Summary: Itemset mining**

- Itemset mining is the simplest of all mining algorithms
- Need to maintain occurrence of each pattern in database
- Tree by lexicographical order is (implicitly) used

# Part 3: Closed Itemset mining



## Solution: Closed Pattern Mining

Find only closed patterns

Observation: Most frequent itemset X can be extended without changing occurrence by adding new elements

def ([Pasquier et al., ICDT'99]). An itemset X is a closed set if and only if there is no proper superset of X with the same frequency (thus the same occurrence set).

### **Closed Pattern Mining**

A closed itemset is the maximal set among all itemsets with the same occurrences.

Equivalence class [X] = {Y| Occ(X)=Occ(Y) }.



## **Brute-force: Stupid Baseline**

#### ALGORITHM Bruteforce

- First, generate all frequent itemsets
- Check them one by one via maximality test
- Maximality test for each candidate frequent set X
  - Add some element e in Σ to X
  - If Freq(X U {e}) is properly less than Freq(X) then reject X.



#### Bruteforce

- STEP1) first, generate all frequent sets
- STEP 2) make closedness test for each set


#### Bruteforce

- STEP1) first, generate all frequent sets
- STEP 2) make closedness test for each set
- STEP3) finally, extract all closed sets



# Complexity of Enumeration Algorithms

Number of patterns usually exponential to Input size input size Input Delay: Time between two pattern outputs Output size Brute-force is exponential delay Delay D w.r.t. pattern size Total Time T

# To achieve linear delay,



# Reverse Search: It's a must

A general mathematical framework to design enumeration algorithms Can be used to prove the correctness of the algorithm Popular in computational geometry Data mining algorithms can be explained in remarkable simplicity

# Often, search space comes as a DAG

- Naive Backtracking
   = Duplication
  - Duplication check by Marking = Exponential Memory
- How to visit all nodes without duplication?



# **Reduction Map**

Mapping from a children to the parent
 Reduction map for closed itemset
 Shrink the itemset until occurrence changes

Take "closure operation"



Closure Operation
 closure(X) of a set X:

 Closed set computed by closure(X) = ∩ { t in T : X ⊆ t }. (taking the intersection of all transactions in T that X occurs as subset)



# Example of Closure Operation

Database T

Non-closed itemset: (B,C)
Occurrence: 1,4
Take Intersection of 1 and 4 (A,B,C,E) ∩ (B,C,E) = (B,C,E)
This is closed itemset



# By applying the reduction map to all nodes, enumeration tree is defined.

 But arrows are in reverse direction..





# **Reverse Search Theorem**

- To prove the correctness, prove the following
  - Reduction map is uniquely defined on all nodes
  - By applying the reduction map repeatedly, one can reach the root node from any node
- Children generation is inverse of reduction map
   Easy to check !

LCM = Linear Time Closed Sets Miner (Uno et al., 2003)

Prefix Preserving **Closure Extension** Jump! = Children generation from the reduction map



# Naïve Closure Extension: Duplication!



#### **Prefix Preserving Closure Extension**

- Ensure any closed set is generated from a unique parent
- **Def.** Closure tail of a closed itemset P $\Leftrightarrow$  the minimum *j* s.t. closure  $(P \cap \{1, ..., j\}) = P$
- **Def.**  $H = \text{closure}(P \cup \{i\})$  is a PPC-extension of *P*   $\Leftrightarrow i > \text{closure tail}$  and  $H \cap \{1, \dots, i-1\} = P \cap \{1, \dots, i-1\}$

#### **Enumeration tree by PPC extension**



# Linear Delay in Pattern Size

(Uno, Uchida, Asai, Arimura, Discovery Science 2004)

<u>Theorem</u>: The algorithm LCM finds all frequent closed sets X appearing in a collection of a transaction database D in O(Imn) time per closed set in the total size of D without duplicates,

where *I* is the maximum length of transactions in **D**, and *n* is the total size of **D**, *m* is the size of pattern **X**.

Note: The output polynomial time complexity of Closed sets discovery is shown by [Makino et al. STACS2002]

# Summary: Closed Itemset Mining

- Closure Extension: Jump from closed set to closed set
- LCM: Linear Delay
- Very fast in practice, too
  - Winner of FIMI'04 (Frequent Itemset Mining Implementation Workshop)
- Relation to clique enumeration (Arimura, Uno, SDM2009)

# Part 4: Ordered Tree Mining

•55

# **Frequent Ordered Tree Mining**

- Natural extension of frequent itemset mining problem for trees
  - Finding all frequent substructure in a given collection of labeled trees
  - How to enumerate them without duplicates

•56

#### Efficient DFS Algorithm

- FREQT [Asai, Arimura, SIAM DM2002]
- TreeMiner [Zaki, ACM KDD2002]
- Rightmost expansion technique

# **Labeled Ordered Trees**

#### Rooted:

#### Ordered:

- Siblings are ordered from left to right.
- Labeled
  - Each node has a label.

#### A model of

- HTML/XML
- Hierarchical records
- Dependency tree of natural language texts.





#### Frequency of a pattern tree

- A root occurrence of pattern T:
  - The node to which the root of T maps by a matching function
- The frequency fr(T) = #root occurrences of T in D



#### **Frequent Tree Mining Problem**

- Given: a colection S of labeled ordered trees and a minimum frequency threshold σ
- Task: Discover all frequent ordered trees in S (with frequency no less than σ) without duplicates







Depth-label sequence in the preorder traversal (depth first search)

 $S = ((d(v_1), I(v_1)), \dots, (d(v_k), I(v_k)))$ 



# **Rightmost expansion**

Extending the DFS Code = Attaching a new node on the rightmost branch
 (d1,l1),...,(dn,ln), (dn+1,ln+1)



#### Searching frequent ordered trees

- Enumerate all frequent ordered trees by backtracking
- Tree extended only by rightmost extension = No duplication



### **Summary: Ordered tree mining**

- Convert tree to a string (DFS Code)
- Adding element to the code = Rightmost extension
- It was relatively easy because nodes are ordered

How about unordered case?

# Part 5: Unordered Tree Mining

# Frequent Unordered Tree Mining

- Unordered trees: Non-trivial subclass of general graphs
- Problem: Exponentially many isomorphic trees
- Efficient DFS Algorithm
  - Unot [Asai, Arimura, DS'03]
  - NK [Nijssen, Kok, MGTS'03]





•68

### **Canonical representation**

# Ordered tree *T* with **lexicographically maximum code**



Left Heavy Condition (Nakano and Uno, 2002)

 T(v): subtree rooted on v
 Ordered tree is canonical *if and only if* Code(T(v1)) ≧Code(T(v2))
 for any pair of sibling nodes v1 (left) and v2 (right)

#### **Reduction Map**

- How to define parent from child in the enumeration tree
- Generate canonical tree of size k-1 from canonical tree of size k
- Remove the last element of DFS Code

Code(T) = (0A, 1B, 2A, 3C, 2B, 1B, 2C)

### **Children Generation**

- Generate children candidates by rightmost extension
- Check the maximality of candidate based on left heavy property
- Discard if not maximal


# Maximality Check by Left Heavy Property

- Code of left subtree must be larger than that of right subtree
- Check only rightmost sibling and second rightmost sibling



## **Complexity of UNOT**

Delay per pattern O(kb<sup>2</sup> mn)
k: pattern size
b: branching factor of the data tree
m: size of data tree
n: database size

# **Summary: Mining Unordered Tree**

The following three elements are necessary for a mining algorithm

- Canonical Representation
- Reduction Map
- Children Generation including Maximality Check

Backtracking on the resulting enumeration tree

# Part 6: Graph Mining

# Frequent Subgraph Mining



• Enumerate all subgraphs occurring more than 3 times



# Gspan (Yan and Han, 2002)

- Most widely used graph mining algorithm
- Can be interpreted with reverse search principle
  - Canonical representation?
  - Reduction map?
  - Children generation?



# DFS Code for Graph

- Depth first search and preorder node labeling
- (src, dest, src\_label, edge\_label, dest\_label)
- Some edges not traversed
  - backward edge (dest < src)</li>
- Elements sorted in the code



# **Canonical Representation**

- Multiple DFS codes: different starting point and children ordering
- Minimum DFS Code: Lexicographically Minimum



((0,1),A,a,A), ((0,2),A,a,A) ((0,1),A,a,A), ((1,2),A,a,A)

# **Reduction Map**

 Removing the tail of minimum DFS code preserves minimality



# **Children Generation**

- Create candidates by adding an element to DFS code
- Check if each candidate is minimum
- If not, remove it



# Minimality Check

- Reconstruct the graph from DFS Code
- Derive the minimum DFS Code by trying all DFSs
  - Speed up by traversing minimal label only
- If the minimal code is not identical to the original, prune it

{(0,1,A,a,A), (1,2,A,a,B), (2,0,B,a,A), (2,3,B,b,A)}



# Summary: Graph Mining

- gSpan is a typical example of reverse search
- Not explained: Closed tree mining, Closed Graph mining
- Delay exponential to pattern size
  - It cannot be avoided due to NP-hardness of graph isomorphism
  - Yet it scales to millions for sparse molecular graphs
- Applications covered in next chapter

# Part7: Dense Module Enumeration

## **Biological Motivation**

- Most cellular processes performed by multi-component protein complexes
- Increasing amount of experimental protein interaction data available
- Our approach
  - Predict complexes (modules) from protein interaction network
  - Exploit additional information given by gene expression data, evolutionary conservation, phenotypic profiles etc.



#### **Protein interaction networks**

- Node: Proteins
- Edge: Physical interaction of two proteins
- Challenge 1: False negative edges
  - Go beyond clique search!
- Challenge 2: False positive edges
  - Assign confidence scores to edges
- Find node sets with high density of high confidence edges

### **Module Discovery**

- Previous work
  - Clique percolation [Palla et al., 2005]
  - Partitioning
    - Hierarchical clustering [Girvan and Newman, 2001]
    - Flow Simulation [Krogan et al., 2006]
    - Spectral methods [Newman, 2006]
  - Heuristic Local Search [Bader and Hogue, 2003]
- Our approach
  - Exhaustively enumerate all dense subgraphs efficiently

# **Motivation for Enumeration Approach**

- Detects overlapping modules
- Allows to specify minimum density for outcoming modules
- Outputs all modules satisfying the density threshold





Partitioning

Enumeration

## **Differential Expression Criterion**

- Incorporation of gene expression
  - Presence of proteins depends on cell type
- Additional Criterion for modules
  - $e_1$ : Num of conditions where whole module expressed
  - $e_0$ : Num of conditions where whole module not expressed
  - Fix minimum values for both quantities



#### **Problem Formalization**

- Interaction network: G = (V, E(V))
- Edge weights:  $0 \le w(\{u, v\}) \le 1$
- Density of  $U \subset V$ :

$$d(U) = \frac{\sum_{\{u,v\}\in E(U)} w(\{u,v\})}{|U|(|U|-1)/2}$$

• Find all  $U \subset V$  with  $d(U) \ge \theta$ ,  $e_1(U) \ge n_1$ , and  $e_0(U) \ge n_0$ 

# **Typical Enumeration Algorithms**

- Itemset mining, graph mining etc.
- Enumerate all entities whose frequency >= 10
- Set up a search tree
- Tree Pruning by anti-monotonicity
  - An entity's frequency is always smaller than that of sub-entity



#### **Network Example**



Density of Modules

- 1,2,3 | 0.5
- 1,3,4 | 0.9
- 2,3,4 | 0.5
- 1,2,4 | 0.3

#### **Graph-shaped Search Space of Modules**



## **Choosing a search tree**

- For efficient search, a search tree is needed
- There are many possible search trees
- Default: Lexicographical ordering



#### Density is not a monotonic criterion

- Subset of dense set is not necessarily dense
- Density does not decrease monotonically on a path
  - Pruning Impossible



#### Question

Is it possible to make a search tree such that density decreases monotonically?

#### Question

Is it possible to make a search tree such that density decreases monotonically?

#### • YES!

Use Reverse Search (Avis and Fukuda 1993)

#### **Reverse search** (Avis and Fukuda, 1993)

- Specify a search tree in the graph-shaped search space
- Reduction Mapping
  - Rule to generate a parent from a child
  - Remove the node with the smallest degree
  - Density always increase by the removal



# Search Tree is uniquely specified by the reduction mapping

 Condition: Every node should converge to the root node by applying the reduction mapping repeatedly



#### Enumeration algorithm by reverse search

- A set of children is generated from a parent node
- Try every possible children, and choose the ones satisfying the reduction mapping
- Prune if no children exist

### **Constraint Integration**

- Differential expression constraint  $e_1(U) \ge n_1$ ,  $e_0(U) \ge n_0$
- Monotonicity: e\_0 and e\_1 decrease with extension of U



Can be used for extra pruning without difficulty

## **Statistical Significance of a module**

- k : The number of nodes in the module
- $\rho$  : Density of the module
- $m_k(\rho)$ : The number of modules of size k with density at least  $\rho$
- Probability of random selection making a denser module (p-value)

$$p = m_k(\rho) / \left( \begin{array}{c} n \\ k \end{array} \right)$$

#### Benchmarking in yeast complex discovery

- Combined interactions from CYGD-Mpact and DIP
- Interactions among 3559 nodes
- Confidence weights on edges due to (Jansen, 2003)
- Methods in comparison
  - Clique detection (Clique)
  - Clique Parcolation Method (CPM)
  - Markov Clustering
- Modules compared with MIPS complexes



## **Evolutionary Conserved Yeast Modules**

- Use ortholog profiles (10 species, InParanoid)
- Density >= 50%, at least three orthologs
- 1917 modules in 30 minutes
- Recovered evolutionary conserved complexes
  - 20S proteasome
  - 19S/22S regulator
  - COPI vesicle coat complex
  - DNA polymerase I and II subunits
  - Translation initiation factor eIF2B complex
- They could not be recovered by simple DME due to low density

# MIPS Complexes discovered by DME (Conserved in Evolution)


#### Phenotype-associated yeast modules

- Use growth phenotypic profiles (21 conditions, Dudley et al, 2005)
- Growth defect in at least one condition
- Each of the 13 highest ranking modules covers the large subunit of mitochondrial ribosome
  - Found additional protein, Mhr1
- Exactly recovered the nucleoplasmic THO complex (Hpr1, Mft1, Rlr1, Thp2)
  - Transcription elongation, hyperrecombination
  - Growth defect under ethanol

109

#### Phenotype Associated Modules Large subunit of mitochondrial ribosome



Mhr1: involved in homologous recombination of the mitochondrial genome

•110

#### **Human Settings**

- Tissue-specific gene expression data (Su et al., 2004)
  - 79 different tissues
- Consistently expressed in 3 tissues, not in 10 tissues
- 7763 proteins, density >= 35%, 5 minutes
- 1021 maximal modules
- MIPS human complex database (Ruepp et al., to appear)

#### Human-expression result

- Around MCM complex, we found inter-complex relationships with ORC, CDC7, Toposome, PLK1 protein
- Module Uqcrc1, Uqcrc2, Uqcrb, Cyc1 (lg p = -13)
  - No overlap with MIPS
  - Ubiquinol-cytochrome c reductase complex
- SCF E3 ubiquitin ligase complex: Mark protein for degradation
  - 5 different modules with different tissue specificity
  - Peripheral proteins: Substrate recognition particles
  - Target proteins are selected in a tissue specific manner!
  - Natural Killer cells have all particles

## High ranking modules around the MCM complex



Expressed in bone mallow cells Not expressed in brain, liver, kidney etc.

## Tissue Specific organization of the SCF ligase complex



#### **Summary: Dense Module Enumeration**

- Novel module enumeration algorithm based on reverse search
- Combination with other information sources
- Statistical significance of dense modules
- Successfully applied to
  - yeast/human protein interaction networks

#### Reference

- [1] R. Agrawal and R. Srikant, "Fast algorithms for mining association rules," in *Proc. 20th Int. Conf. Very Large Data Bases*, 1994.
- [2] T. Asai, H. Arimura, T. Uno, and S. Nakano, "Discovering frequent substructures in large unordered trees," *Discovery science*, 2003.
- [3] T. Asai, K. Abe, S. Kawasoe, H. Arimura, H. Sakamoto, and S. Arikawa, "Efficient Substructure Discovery from Large Semi-structured Data.," in SDM, 2002.
- [4] D. Avis and K. Fukuda, "Reverse search for enumeration," Discrete Applied Mathematics, vol. 65, no. 1-3, pp. 21-46, 1996.
- [5] E. Georgii, S. Dietmann, T. Uno, P. Pagel, and K. Tsuda, "Enumeration of conditiondependent dense modules in protein interaction networks.," *Bioinformatics (Oxford, England)*, vol. 25, no. 7, pp. 933-40, Apr. 2009.
- [6] S. Nijssen and J. Kok, "Efficient discovery of frequent unordered trees," *MGTS*, 2003.
- [7] N. Pasquier, Y. Bastide, R. Taouil, and L. Lakhal, "Discovering frequent closed itemsets for association rules," *Database Theory–ICDT*'99, 1999.
- [8] I. Takigawa and H. Mamitsuka, "Efficiently mining δ-tolerance closed frequent subgraphs," *Machine Learning*, vol. 82, no. 2, pp. 95-121, Sep. 2010.
- [9] X. Yan and J. Han, "gSpan: graph-based substructure pattern mining," in 2002 IEEE International Conference on Data Mining, 2002. Proceedings., pp. 721-724.
- [10] M. Zaki, "Efficiently mining frequent trees in a forest: Algorithms and applications," IEEE Transaction on Knowledge and Data Engineering, 2005.

•116

Machine Learning Summer School Kyoto 2012 Graph Mining Chapter 2: Learning from Structured Data

National Institute of Advanced Industrial Science and Technology

Koji Tsuda

## Agenda 2 (Learning from Structured Data)

1. Preliminaries
2. Graph Clustering by EM
3. Graph Boosting
4. Regularization Paths in Graph Classification
5. Itemset Boosting for predicting HIV drug resistance

## Part 1: Preliminaries

## **Clustering Graphs**



















































## Substructure Representation

0/1 vector of *pattern* indicators
Huge dimensionality!
Need Graph Mining for selecting features



## **Graph Mining**

## Frequent Substructure Mining Enumerate all patterns occurred in at least m graphs

 $n_{-}$ 

$$S_{freq} = \{k \mid \sum_{i=1}^{n} x_{ik} \ge m\}.$$

 $x_{ik} \in \{0,1\}$  :Indicator of pattern k in graph i

Support(k): # of occurrence of pattern k

Enumeration on Tree-shaped Search Space

- Each node has a pattern
- Generate nodes from the root:
  - Add an edge at each step



Support(g): # of occurrence of pattern g

## **Tree Pruning**

Anti-monotonicity:

$$g \subseteq g' \implies support(g) \ge support(g')$$

If support(g) < m, stop exploring!</pre>



#### Gspan (Yan and Han, 2002)





Efficient detection of isomorphic patterns



## Depth First Search (DFS) Code



## **Discriminative patterns**

w\_i > 0: positive class
 w\_i < 0: negative class</li>
 Weighted Substructure Mining

$$S_w = \{k \mid \left| \sum_{i=1}^n w_i (2x_{ik} - 1) \right| \ge \tau \},\$$

Patterns with large frequency difference
 Not Anti-Monotonic: Use a bound

## **Multiclass version**

# Multiple weight vectors $w_{\ell i} > 0$ (graph *i* belongs to class $\ell$ ) $w_{\ell i} < 0$ (otherwise)

Search patterns overrepresented in a class

$$S_W = \{k \mid \max_{\ell=1,\dots,c} \left| \sum_{i=1}^n w_{\ell i} (2x_{ik} - 1) \right| - \tau_{\ell} \ge 0 \}.$$

## **Basic Bound**

- $x_{ij}$ : Occurrence of pattern j
- If k is supergraph of pattern j,

$$\begin{aligned} \sum_{i=1}^{n} w_{\ell i}(2x_{ik} - 1) &| \leq \gamma_{\ell} \qquad \gamma_{\ell} = \max(\gamma_{\ell}^{+}, \gamma_{\ell}^{-}) \\ \gamma_{\ell}^{+} &= 2 \sum_{\{i \mid w_{\ell i} \geq 0, x_{ij} = 1\}} |w_{\ell i}| - \sum_{i=1}^{n} w_{\ell i} \\ \gamma_{\ell}^{-} &= 2 \sum_{\{i \mid w_{\ell i} < 0, x_{ij} = 1\}} |w_{\ell i}| + \sum_{i=1}^{n} w_{\ell i} \end{aligned}$$

## **Pruning Condition**

| n

#### Summarizing the bound for all classes,

$$\max_{\ell} \left| \sum_{i=1}^{n} w_{\ell i} (2x_{ik} - 1) \right| - \tau_{\ell} \le \max_{\ell} (\gamma_{\ell} - \tau_{\ell})$$

If it is negative, the search tree can be pruned safely

## Summary: Preliminaries

Various graph learning problems
 Supervised/Unsupervised

Discovery of salient features by graph mining
 Actual speed depends on the data
 Faster for..
 Sparse graphs
 Diverse kinds of labels

# Part 2: EM-based clustering of graphs

## EM-based graph clustering

#### Motivation

- Learning a mixture model in the feature space of patterns
- Basis for more complex probabilistic inference
- L1 regularization & Graph Mining
   E-step -> Mining -> M-step

## **Probabilistic Model**

Binomial Mixture

$$p(\boldsymbol{x}|\Theta) = \sum_{\ell=1}^{c} \alpha_{\ell} p_{\ell}(\boldsymbol{x}|\boldsymbol{\theta}_{\ell})$$

Each Component

$$p_{\ell}(\boldsymbol{x}|\boldsymbol{\theta}_{\ell}) = \prod_{k=1}^{d} \frac{\exp(\theta_{\ell k} x_{k})}{1 + \exp(\theta_{\ell k})}$$

 $oldsymbol{x}$  :Feature vector of a graph (0 or 1)  $lpha_\ell$  :Mixing weight for cluster  $\ell$  $oldsymbol{ heta}_\ell$  :Parameter vector for cluster  $\ell$ 

## **Ordinary EM algorithm**





◆ E-step: Get posterior r<sub>ℓi</sub> = p(y = ℓ|x<sub>i</sub>)
◆ M-step: Estimate θ<sub>ℓ</sub> using posterior probs.
◆ Both are computationally prohibitive (!)

## Regularization



#### $\bullet$ Baseline constant $oldsymbol{ heta}_0$

ML parameter estimate using single binomial distribution

$$\theta_{0k} = \log \eta_{0k} - \log(1 - \theta_{0k})$$
  $\eta_{0k} = \frac{1}{n} \sum x_{ik}$ 

In solution, most parameters exactly equal to constants

#### E-step

#### Active pattern

 $F = \{k \mid \text{there exists } \ell \text{ such that } \theta_{\ell k} \neq \theta_{0k} \}.$ 

# E-step computed only with active patterns (computable!)

$$p(y = \ell | \boldsymbol{x}) = \frac{\alpha_{\ell} \prod_{k \in F} p_{\ell k}(x_k | \theta_{\ell k})}{\sum_{\ell} \alpha_{\ell} \prod_{k \in F} p_{\ell k}(x_k | \theta_{\ell k})}$$

#### M-step

- Putative cluster assignment  $r_{\ell i}$
- Each parameter is solved separately

$$\min_{\theta_{\ell k}} -\frac{1}{n} \sum_{i} r_{\ell i} \log p(x_{ik} | \theta_{\ell k}) + \lambda |\theta_{\ell k} - \theta_{0k}|.$$

- Naïve way:
  - solve it for all params and identify active patterns

#### Use graph mining to find active patterns

### Solution

$$\theta_{\ell k} = \begin{cases} \log \frac{\eta_{\ell k} - \lambda_{\ell}}{1 - (\eta_{\ell k} - \lambda_{\ell})} & (\eta_{\ell k} \ge \eta_{0k} + \lambda_{\ell}). \\ \theta_{0k} & (\eta_{0k} - \lambda_{\ell} \le \eta_{\ell k} \le \eta_{0k} + \lambda_{\ell}) \\ \log \frac{\eta_{\ell k} + \lambda_{\ell}}{1 - (\eta_{\ell k} + \lambda_{\ell})} & (\eta_{\ell k} \le \eta_{0k} - \lambda_{\ell}). \end{cases}$$

Occurrence probability in a cluster

$$\eta_{\ell k} = \sum r_{\ell i} x_{ik} / \sum r_{\ell j}$$

Overall occurrence probability

l

$$\eta_{0k} = \frac{1}{n} \sum_{i} x_{ik}$$

 $\lambda_{\ell} = \frac{\lambda n}{\sum_{j} r_{\ell j}}$ 

#### Solution



## **Important Observation**

For active pattern k, the occurrence probability in a graph cluster is significantly different from the average

$$\theta_{\ell k} \neq \theta_{0k} \Leftrightarrow |\eta_{\ell k} - \eta_{0k}| \ge \lambda$$



26

## Mining for Active Patterns

Active pattern  $F = \{k \mid \text{there exists } \ell \text{ such that } \theta_{\ell k} \neq \theta_{0k} \}.$ Equivalently written as  $F = \{k \mid \max_{\ell=1,\dots,c} \left| \sum_{i=1}^{n} w_{\ell i} (2x_{ik} - 1) \right| - 2\lambda_{\ell} \ge 0 \}.$ F can be found by graph mining! (multiclass)

## Experiments: RNA graphs

Stem as a node
Secondary structure by RNAfold
0/1 Vertex label (self loop or not)




## **Clustering RNA graphs**

Three Rfam families Intron GP I (Int, 30 graphs) SSU rRNA 5 (SSU, 50 graphs) RNase bact a (RNase, 50 graphs) Three bipartition problems Results evaluated by ROC scores (Area under the ROC curve)

## **Examples of RNA Graphs**



## **ROC Scores**

н

	Int-SSU	Int-RNase	SSU-RNase
MGK	0.748	0.531	0.878
Spec	0.550	0.573	0.848
$\lambda = 0.01$	0.824	0.921	0.863
$\lambda = 0.02$	0.821	0.920	0.862
$\lambda = 0.03$	0.825	0.948	0.843
$\lambda = 0.04$	0.832	0.947	0.825
$\lambda = 0.06$	0.831	0.941	0.782
$\lambda = 0.08$	0.845	0.941	0.787
$\lambda = 0.10$	0.815	0.927	0.786

#### No of Patterns & Time

ı

	Int-SSU	Int-RNase	SSU-RNase
$\lambda = 0.01$	12505 (71s)	14366 (77s)	$17934 \ (102s)$
$\lambda = 0.02$	12596~(75s)	10988~(65s)	11025 (76s)
$\lambda = 0.03$	9799~(66s)	7632~(52s)	8875 (73s)
$\lambda = 0.04$	6904 (57s)	5924 (45s)	6925~(67s)
$\lambda = 0.06$	5093 (47s)	4305~(37s)	5230 (58s)
$\lambda = 0.08$	4065 (42s)	3001 (32s)	3896~(50s)
$\lambda = 0.10$	3245 (37s)	2074 (26s)	2923 (44s)

### **Found Patterns**









1.117

1.117

















33

## Summary (graph EM)

- Substructure representation is better than paths
- Probabilistic inference helped by graph mining
- Extension to Dirichlet mixture model
  - Reported in Tsuda et al., SDM 2008
- Possible extension
  - Graph PCA, LFD, CCA
  - Semi-supervised learning

#### Part 3: Graph Boosting

## Graph classification problem in chemoinformatics

- Known as SAR problem in chemical informatics
  - Quantitative) Structure-Activity Analysis
- Given a graph, predict a class-label (+1 or -1)
   Typically, features (descriptors) are given
   e.g., Dragon Descriptors, JOELIB2

#### SAR with conventional descriptors

	#atoms	#bonds	#rings	 Activity
HO CI CH, OH	22	25		+1
	20	21		+1
HO CHI CON	23	24		+1
H <sub>i</sub> C + CH <sub>i</sub>	11	11		-1
S C C C C C C C C C C C C C C C C C C C	21	22		-1

## **Motivation of Graph Boosting**

- Descriptors are not always available
- New features by obtaining informative patterns (i.e., subgraphs)
- Greedy pattern discovery by Boosting + gSpan
- Linear Programming (LP) Boosting
  - Reduce the number of graph mining calls
  - Faster than AdaBoost
- Accurate prediction & interpretable results



#### SAR with patterns



Sparse classification in a very high dimensional space

 G: all possible patterns (intractably large) |G|-dimensional feature vector x for a molecule

Linear Classifier  $f(\mathbf{x}) = \sum_{j=1}^{d} \alpha_j x_j$ 



Use L1 regularizer to have sparse a Select a tractable number of patterns

#### **Problem formulation**



Sum of hinge loss and L1 regularizer

- $\{m{x}_n, y_n\}_{n=1}^\ell$ : Training examples
- ξ: slack variables

#### Dual LP

Primal: Huge number of weight variablesDual: Huge number of constraints

Dual problem  $\begin{array}{ll}
\max_{\mathbf{u}} & \sum_{n=1}^{\ell} u_i \\
s.t. & \sum_{n=1}^{\ell} u_n y_n x_{ni} \leq 1, \ i = 1, \dots, d \\
& 0 \leq u_n \leq C, \ n = 1, \dots, \ell
\end{array}$  Column Generation Algorithm for LP Boost (Demiriz et al., 2002)

Start from the dual with no constraints
 Add the most violated constraint each time
 Guaranteed to converge



# Finding the most violated constraint

Constraint for a pattern (shown again)

$$\sum_{n=1}^{\ell} u_n y_n x_{ni} \le$$

Finding the most violated one

0

$$\operatorname{argmax}_i \sum_{n=1}^{\ell} u_n y_n x_{ni}$$

Searched by weighted substructure mining

## **Algorithm Overview**

#### Iteration

- Find a new pattern by graph mining with weight **u**
- If all constraints are satisfied, break
- Add a new constraint
  - Update **u** by solving the dual problem
- Return
  - Convert dual solution to obtain primal solution a

#### **Experimental Settings**

#### Classification and Regression Datasets

	# da	ita #	<sup>∉</sup> positives	# negative	es avg. atoms	avg. bonds
 CPDB	<b>68</b> 4	1	341	343	14.1	14.6
CAS	433	7	2401	1936	29.9	30.9
	•		#data	avg. atoms	avg. bonds	
		AR	146	19.5	21.1	
		$\mathbf{ER}$	131	19.2	20.7	
		$\mathbf{ES}$	59	18.2	19.7	
			- 			

#### **Classification Results**

Table 2 Classification performance obtained by 10-fold cross validation in two datasets measured by the accuracy (ACC) and the area under the ROC curve (AUC). We obtained the results of MGK and gBoost from our implementations, but the other results are quoted from literature. The best results are highlighted in bold fonts.

	Measure	MOLFEA[13]	Gaston[19]	CPM[2]	MGK[18]	gBoost
CPDB	ACC	-	78	75.96	76.5	78.8
	AUC	-	-	-	0.756	0.854
CAS	ACC	79	-	80.14	77.1	82.5
	AUC	-	-	-	0.763	0.889
1 1 1	-	- 5 1 5 5 5				
						48

#### **Regression Results**

Table 3 Regression performance obtained by leave-one-out cross validation in three assays from the EDKB evaluated by mean absolute error (MAE), root mean squared error (RMSE), and  $Q^2$ . Note that for MAE and RMSE, lower values indicate better prediction, which is vice versa for  $Q^2$ . We obtained the results of MGK and gBoost from our implementations, but the other results are quoted from literature. The best results are highlighted in bold fonts.

	Measure	CoMFA[14, 37]	MGK[18]	gBoost
AR	MAE	-	0.229	0.176
	RMSE	-	0.335	0.232
	$Q^2$	0.571	0.346	0.682
$\mathbf{ER}$	MAE	-	0.320	0.307
	RMSE	-	0.427	0.393
	$Q^2$	0.660	0.267	0.378
$\mathbf{ES}$	MAE	-	0.322	0.249
	RMSE	-	0.413	0.362
	$Q^2$	-	0.522	0.632

#### Extracted patterns from CPDB



50

#### Memory Usage



#### Runtime



52

#### **Comparison with AdaBoost**



## Summary (Graph Boosting)

 Graph Boosting simultaneously generate patterns and learn their weights
 Finite convergence by column generation
 Interpretable by chemists.
 Flexible constraints and speed-up by LP.

## Part 4: Entire Regularization

#### Path

#### Overview

#### Entire regularization paths LARS-LASSO (Efron et al., 2004), L1SVM Forward selection of features Trace the solution trajectory of L1-regularized learning Path following algorithm for graph data Feature search -> pattern search Branch-and-bound algorithm DFS code tree, New Bound

## Path Following Algorithms

• LASSO regression  $\beta(\lambda) = \underset{\beta}{\operatorname{argmin}} L(\boldsymbol{y}, X\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1.$ • Follow the complete trajectory of  $\boldsymbol{\beta}(\lambda)$ •  $\lambda$  : Infinity to Zero

#### $\bullet$ Active feature set $\mathcal{A}$

Features corresponding to nonzero weights

#### **Piecewise Linear Path**

At a turning point,

or

- A new feature included into the active set,
- An existing feature excluded from the active set



#### Practical Merit of Path Following

- Cross validation by grid search
   Has to solve QP many times
   Especially time-consuming for graph data
- Path following does not include QP
   Determine the CV-optimal regularization parameter in the finest precision

#### Pseudo code of path following



#### Feature space of patterns

♦ Graph training data G = {G<sub>i</sub>}<sup>n</sup><sub>i=1</sub>
♦ Set of all subgraphs (patterns) T
♦ Each graph is represented as

$$\boldsymbol{x}_i = (x_{it})_{t \in T}, x_{it} = I(t \subseteq G_i)$$

$$\begin{array}{c} \mathbb{B} \\ (0, \ldots, 0, 1, 0, \ldots, 0, 1, 0, \ldots) \\ \hline \\ \mathbb{A} \\ \hline \\ \mathbb{A} \\ \mathbb{B} \\ \mathbb{A} \end{array}$$

#### Main Search problem

#### Step size if pattern t is included next

$$d_t = \min_{+} \left\{ \frac{\rho_0 - \sum_i w_i x_{it}}{\eta_0 - \sum_i v_i x_{it}}, \quad \frac{\rho_0 + \sum_i w_i x_{it}}{\eta_0 + \sum_i v_i x_{it}} \right\}$$

 $w_i, v_i, \rho_0, \eta_0$ : constants computed from active set

Find pattern  $t \in \mathcal{T}$  that minimizes  $d_t$ 

#### **Tree-shaped Search Space**

- Each node has a pattern
   Generate nodes from the root:
  - Add an edge at each step



#### **Tree Pruning**

If it is guaranteed that the optimal pattern is not in the downstream, the search tree can be pruned


#### Theorem (Pruning condition)





#### No better pattern in the downstream, if

$$b_{w} + d_{t}^{*}b_{v} < |\rho_{0}| - d_{t}^{*}|\eta_{0}|$$

where

$$b_{\boldsymbol{w}} = \max\left\{\sum_{w_i < 0} |w_i| x_{it}, \sum_{w_i > 0} |w_i| x_{it}\right\}$$

#### **Initial Experiments**

EDKB Estrogen receptor database
 131 training graphs (chemical compounds)

Computation Time: 4 sec/search
 Pattern size limitation: 10 edges

#### **Solution Path**





Summary: Regularization Path Path following implemented for graph data Pattern search by graph mining Classification: To do Combination with item set mining

# Part 5: . Itemset Boosting for predicting HIV drug resistance

#### Life cycle of an HIV Virus



## Approved HIV Drugs

- 8 Protease inhibitors (PI)
- 8 Nucleoside/nucleotide reverse transcriptase inhibitors (NRTI)
- 3 Non-nucleoside reverse transcriptase inhibitors (NNRTI)
- 1 Fusion inhibitor

#### Drug resistance of HIV

- Exposure to a drug causes mutations in HIV's genes
- As a result, HIV gains resistance against the drug
- Cost of identifying the genotypes of HIV in a patient is relatively cheap
- Predict the drug resistance from HIV's genotypes !
   Effective Pharmacotherapy for individuals

## Drug resistance prediction problem as regression

- Input: Set of mutations in a target protein

   41L: Amino acid at position 41 changed to L
- Output: Resistance against a drug (fold change)

(40F,41L,43E,210W,211K,215Y) ⇒ 0.8 (43Q,75M,122E,210W,211K) ⇒ 12.8 (75I, 77L, 116Y,151M,184V) ⇒ ?

#### Simple vs Complex Genotypic Features

• Simple genotypic features

Complex genotypic features

#### Linear Regression on Simple Genotypic Features (Rhee et al., PNAS 2006)

Mutation associations not discovered



#### **Complex Genotypic Features**

- 0/1 vector of pattern indicators
- Huge dimensionality!
- Need itemset mining for selecting features
- Selection of salient features



77L,116Y 103N,210W,215Y

patterns

#### Other methods

- Nonlinear SVM
   Not interpretable
  - High accuracy
- Decision trees
  - Interpretable
  - Low accuracy



 $\begin{aligned} x_1 \wedge x_6 \Rightarrow y_1, \ x_1 \wedge \neg x_6 \Rightarrow y_2, \ \neg x_1 \wedge \neg x_7 \Rightarrow y_5 \\ \neg x_1 \wedge x_7 \wedge x_2 \Rightarrow y_3, \ \neg x_1 \wedge x_7 \wedge \neg x_2 \Rightarrow y_4. \end{aligned}$ 

#### **Motivation of Itemset Boosting**

- Impossible to maintain all complex features
   Greedy feature discovery by Boosting + itemset mining
- Quadratic Programming (QP) Boosting
  - Reduce the number of itemset mining calls
  - Faster than AdaBoost
- Accurate prediction & interpretable results

Sparse classification in a very high dimensional space

 G: all possible patterns (intractably large) |G|-dimensional feature vector **x** 

Linear Classifier

 $f(\mathbf{x}) = \sum_{i=1}^{a} \alpha_{i} x_{j}$ 



Use L1 regularizer to have sparse a (LASSO)



Select a tractable number of patterns

#### Problem formulation: Quadratic Programming

$$\min_{\boldsymbol{\alpha},\boldsymbol{\xi}} \quad ||\boldsymbol{\alpha}||_1 + \frac{C}{2} \sum_{n=1}^{\ell} \xi_n^2$$

s.t.  $|y_n - \boldsymbol{\alpha}^\top \mathbf{x}_n| \le \xi_n, \quad \xi_n \ge 0, \quad n = 1, \dots, \ell$ 

Sum of squared loss and L1 regularizer

- $\{m{x}_n, y_n\}_{n=1}^\ell$ : Training examples
- ξ: slack variables

## Dual QP

Primal: Huge number of weight variablesDual: Huge number of constraints

$$\min_{\mathbf{u}} \qquad \frac{1}{2C} \sum_{n=1}^{\ell} u_n^2 - \sum_{n=1}^{\ell} y_n u_n \\ s.t. \qquad -1 \le \sum_{n=1}^{\ell} u_n x_{ni} \le 1, \ i = 1, \dots, d \\ \sum_{n=1}^{\ell} u_i = 0$$

Column Generation Algorithm for QP Boost (Demiriz et al., 2002)

Start from the dual with no constraints
 Add the most violated constraint each time
 Guaranteed to converge

Used Part Finding the most violated constraint

#### Constraint for a pattern (shown again)

 $\left|\sum_{n=1}^{\ell} u_n x_{ni}\right| \le 1$ 

Finding the most violated one

$$\operatorname{argmax}_{i} \left| \sum_{n=1}^{\ell} u_{n} x_{ni} \right|$$

Searched by weighted itemset mining

#### **Algorithm Overview**

#### Iteration

- Find a new pattern by graph mining with weight **u**
- If all constraints are satisfied, break
- Add a new constraint
  - Update **u** by solving the dual problem
- Return
  - Convert dual solution to obtain primal solution a

### Experimental settings

- Three classes of drugs
  - NRTI (Nucleotide Reverse Transcriptase Inhibitors)
  - PI (Protease Inhibitors)
  - NNRTI (Nonnucleotide Reverse Transcriptase Inhibitors)
- 5fold cross validation
  - Linear SVM, Ridge regression, Lars
  - Nonlinear SVM, iBoost

#### **Regression results**

Drug	# isolates	Linear Methods			Nonlinear Methods					
		SVR linear	Ridge	LARS	SVR p2	SVR p3	SVR r1	SVR r10	SVR r100	iBoost
NRTI										
Lamivudine (3TC)	633	0.913	0.753	0.93	0.927	0.876	0.306	0.934	0.608	0.940
Abacavir (ABC)	628	0.731	0.585	0.79	0.772	0.720	0.256	0.745	0.614	0.801
Zidovudine (AZT)	630	0.751	0.72	0.65	0.797	0.760	0.240	0.754	0.515	0.797
Stavudine (D4T)	630	0.729	0.662	0.76	0.729	0.676	0.184	0.732	0.534	0.789
Didanosine (DDI)	632	0.695	0.591	0.72	0.690	0.653	0.170	0.705	0.372	0.737
Tenofovir (TDF)	353	0.603	0.568	0.40	0.554	0.465	0.107	0.612	0.366	0.552
Average		0.737	0.647	0.708	0.744	0.692	0.211	0.747	0.502	0.769
NNRTI										
Delavirdine (DLV)	732	0.799	0.794	0.79	0.729	0.684	0.189	0.802	0.323	0.771
Efavirenz (EFV)	734	0.793	0.772	0.85	0.730	0.640	0.170	0.797	0.205	0.771
Nevirapine (NVP)	746	0.757	0.719	0.79	0.704	0.592	0.166	0.765	0.181	0.781
Average		0.783	0.762	0.810	0.721	0.639	0.175	0.788	0.236	0.774
PI										
Amprenavir (APV)	768	0.819	0.749	0.81	0.756	0.697	0.483	0.82	0.576	0.802
Atazanavir (ATV)	329	0.724	0.594	0.76	0.731	0.654	0.297	0.739	0.450	0.701
Indinavir (IDV)	827	0.830	0.710	0.81	0.819	0.775	0.574	0.844	0.710	0.816
Lopinavir (LPV)	517	0.845	0.715	0.85	0.821	0.757	0.496	0.857	0.697	0.827
Nelfinavir (NFV)	844	0.845	0.697	0.84	0.800	0.728	0.615	0.858	0.638	0.838
Ritonavir (RTV)	795	0.893	0.794	0.89	0.875	0.820	0.598	0.903	0.788	0.892
Saquinavir (SQV)	826	0.814	0.738	0.81	0.782	0.733	0.468	0.829	0.559	0.791
Average		0.824	0.712	0.824	0.798	0.738	0.504	0.836	0.631	0.810

Accuracy: 
$$r^2 = 1 - \frac{\sum_{i=1}^n (y_i - f(x_i))^2}{\sum_{i=1}^n (y_i - \frac{1}{n} \sum_{i=1}^n y_i)^2}$$
. 87

## Accuracy Summary

- NRTIs: iBoost performed best
- Pls:
  - Nonlinear methods were better than linear
  - SVMs were slightly better (non-significant)
- NNRTIs
  - Linear methods were better
  - Combination is not necessary



# Known mutation associations in RT

- Red: Thymidineassociated Mutations (TAM)
  - 41L, 67N, 70R, 210W,
    215Y/F, 219Q, 69i
- Blue: Q151M Complex
  - 75I, 77L, 116Y, 151M,
    65R, 74V, 184I/V



RT of HIV-1 90





## NNRTI Drugs (almost no combination found)

## **Computation Time of iBoost**

- Training time for 3TC
- 507 isolates with 371 mutations on average
- QP time longer than mining time



## Summary (HIV)

- Itemset Boosting for finding mutation associations
- Good accuracy for NRTIs
- Our complex features re-discover known mutation clusters
- Broad applications
  - Multiple SNP analysis, RNAi efficacy prediction
  - Motif combination, Flu mutation analysis
  - P53 mutation analysis

#### Reference

- [1] K. Tsuda and T. Kudo, "Clustering graphs by weighted substructure mining," in *Proceedings of the 23rd international conference on Machine learning ICML '06*, 2006, pp. 953–960.
- [2] H. Saigo, N. Krämer, and K. Tsuda, "Partial least squares regression for graph mining," in *Proceeding of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining KDD '08*, 2008, p. 578.
- [3] H. Saigo, S. Nowozin, T. Kadowaki, T. Kudo, and K. Tsuda, "gBoost: a mathematical programming approach to graph classification and regression," *Machine Learning*, vol. 75, no. 1, pp. 69–89, Nov. 2008.
- [4] H. Saigo, T. Uno, and K. Tsuda, "Mining complex genotypic features for predicting HIV-1 drug resistance.," *Bioinformatics (Oxford, England)*, vol. 23, no. 18, pp. 2455–62, Sep. 2007.
- [5] K. Tsuda, "Entire regularization paths for graph data," in *Proceedings of the 24th international conference on Machine learning ICML '07*, 2007, pp. 919–926.

Machine Learning Summer School Kyoto 2012 Graph Mining Chapter 3: Kernels for Structured Data

National Institute of Advanced Industrial Science and Technology

Koji Tsuda

## Agenda 3 (Kernel)

1. Kernel Method Revisited
2. Marginalized Kernels (Fisher Kernels)
3. Marginalized Graph Kernels
4. Weisfeiler-Lehman kernels
5. Reaction Graph kernels
6. Concluding Remark

#### Part 1: Kernel Method Revisited

#### **Kernels and Learning**

 In Kernel-based learning algorithms, problem solving is now decoupled into:
 A general purpose learning algorithm (e.g. SVM, PCA, ...) – Often linear algorithm
 A problem specific kernel

Complex Learning Task Simple (linear) learning algorithm

Specific Kernel function

#### **Current Synthesis**

#### Modularity and re-usability

- Same kernel ,different learning algorithms
- Different kernels, same learning algorithms


# Kernel Methods : intuitive idea

- Kernel represents the similarity between two objects, defined as the dot-product in this new vector space
- But the mapping is left implicit
- Easy generalization of a lot of dotproduct-based learning algorithms

# Kernel Methods : the mapping



# Kernel : more formal definition

## A kernel k(x,y)

- is a similarity measure
- defined by an implicit mapping  $\phi$ ,
- such that:  $k(x,y) = \phi(x) \cdot \phi(y)$

### This similarity measure implies:

- Invariance or other a priori knowledge
- The class of functions the solution is taken from
  - Possibly infinite dimension (hypothesis space for learning)

 ... but still computational efficiency when computing k(x,y)

# Kernel Trick

- Generalizes (nonlinearly) algorithms in clustering, classification, density estimation ...
  - When these algorithms are dot-product based, by replacing the dot product (x·y) by k(x,y)=\u03c6(x)·\u03c6(y)
  - When these algorithms are distance-based, by replacing d(x,y) by k(x,x)+k(y,y)-2k(x,y)

# Valid Kernels

Theorem: k(x,y) is a valid kernel if k is positive definite and symmetric (Mercer Kernel)

- A function is P.D. if  $\int K(\mathbf{x}, \mathbf{y}) f(\mathbf{x}) f(\mathbf{y}) d\mathbf{x} d\mathbf{y} \ge 0 \quad \forall f \in L_2$
- In other words, the Gram matrix K (whose elements are k(x<sub>i</sub>,x<sub>j</sub>)) must be positive definite for all x<sub>i</sub>, x<sub>j</sub> of the input space
- One possible choice of  $\phi(x)$ :  $k(\cdot,x)$  (maps a point x to a function  $k(\cdot,x) \rightarrow$  feature space with infinite dimension!)

# How to build new kernels Kernel combinations, preserving validity: $K(\mathbf{x},\mathbf{y}) = \lambda K_1(\mathbf{x},\mathbf{y}) + (1-\lambda)K_2(\mathbf{x},\mathbf{y}) \quad 0 \le \lambda \le 1$ $K(\mathbf{x}, \mathbf{y}) = a.K_1(\mathbf{x}, \mathbf{y}) \quad a > 0$ $K(\mathbf{x},\mathbf{y}) = K_1(\mathbf{x},\mathbf{y}).K_2(\mathbf{x},\mathbf{y})$ $K(\mathbf{x}, \mathbf{y}) = f(x) \cdot f(y)$ f is real-valued function $K(\mathbf{x},\mathbf{y}) = K_3(\boldsymbol{\varphi}(\mathbf{x}),\boldsymbol{\varphi}(\mathbf{y}))$ $K(\mathbf{x}, \mathbf{y}) = \mathbf{x}' P \mathbf{y}$ P symmetric definite positive $K(\mathbf{x}, \mathbf{y}) = \frac{K_1(\mathbf{x}, \mathbf{y})}{\sqrt{K_1(\mathbf{x}, \mathbf{x})}\sqrt{K_1(\mathbf{y}, \mathbf{y})}}$

# Strategies of Design

Convolution Kernels: text is a recursively-defined data structure. How to build "global" kernels form local (atomic level) kernels?

Generative model-based kernels: the "topology" of the problem will be translated into a kernel function

### Family of kernels

### Kernels for biological sequences

- Spectrum kernel
- Marginalized kernel
- Profile kernel
- Local alignment kernel

#### Tree Kernels

- Kernel for phylogenetic profiles
- Kernel for natural language
- Kernel for RNA sequences



Kernel Methods in Computational

**Biology, MIT Press** 

### Family of kernels

### Kernels for nodes in a network

- Diffusion kernel
- Locally constrained diffusion kernel

#### Graph Kernels

- Marginalized Graph Kernels
- MGK without tottering
- Acyclic Pattern Kernels
- Shortest Path Kernel
- Weisfeiler-Lehman Kernel

## Weak points of kernel methods

Not Interpretable
 Not sure which features are used
 -> Graph Mining, Boosting

Dense kernel matrices: Slow
 Take O(n^3) time for manipulation
 -> Semi-supervised learning

# Part 2. Marginalized kernels

# Biological Sequences: Classification Tasks

DNA sequences (A,C,G,T)

 Gene Finding, Splice Sites

 RNA sequences (A,C,G,U)

 MicroRNA discovery, Classification into Rfam families

 Amino Acid Sequences (20 symbols)

 Remote Homolog Detection, Fold

recognition

# **Kernels for Sequences**

Similarity between sequences of different lengths

## ACGGTTCAA

## ATATCGCGGGAA

Later combined with SVMs and other kernel methods

# Count Kernel

Inner product between symbol counts

G T A 3 2 2 2ACGGTTCAA ATATCGCGGGAA 4 2 4 2

Extension: Spectrum kernels (Leslie et al., 2002)

Counts the number of k-mers (k-grams) efficiently

 $\mathbf{C}$ 

Not good for sequences with <u>frequent context change</u> E.g., coding/non-coding regions in DNA

## Hidden Markov Models for Estimating Context

 ${}$  Visible Variable  $oldsymbol{x} = (x_1, \ldots, x_m)$  : Symbol Sequence

• Hidden Variable  $h = (h_1, \ldots, h_m)$ : Context

• HMM can estimate the posterior probability of hidden variables p(h|x) from data

# h: 1 2 2 1 2 2 1 2 2 x: A C G G T T C A A

# Marginalized kernels

Design a joint kernel K<sub>z</sub>(z, z') for combined z = (x, h)
 Hidden variable is not usually available

 Take expectation with respect to the hidden variable

The marginalized kernel for visible variables

$$K(\boldsymbol{x}, \boldsymbol{x}') = \sum_{\boldsymbol{h} \in \mathcal{H}} \sum_{\boldsymbol{h}' \in \mathcal{H}} p(\boldsymbol{h} | \boldsymbol{x}) p(\boldsymbol{h}' | \boldsymbol{x}') K_{\boldsymbol{z}}(\boldsymbol{z}, \boldsymbol{z}')$$

# Designing a joint kernel for sequences

Symbols are counted separately in each context

h:122122122x:ACGGTTCA(A,1) = 1(C,1) = 1(G,1) = 1(T,1) = 0x:ACGGTTCA(A,2) = 2(C,2) = 1(G,2) = 1(T,2) = 2

•  $c_{k\ell}(\boldsymbol{z})$  :count of a combined symbol (k,l)

Joint kernel: count kernel with context information

$$K_{z}(\boldsymbol{z}, \boldsymbol{z}') = \sum_{k=1}^{n_{x}} \sum_{\ell=1}^{n_{h}} c_{k\ell}(\boldsymbol{z}) c_{k\ell}(\boldsymbol{z}')$$

## Marginalization of the joint kernel

• Joint kernel
$$K_z(\boldsymbol{z}, \boldsymbol{z}') = \sum_{k=1}^{n_x} \sum_{\ell=1}^{n_h} c_{k\ell}(\boldsymbol{z}) c_{k\ell}(\boldsymbol{z}')$$

Marginalized count kernel

$$K(\boldsymbol{x}, \boldsymbol{x}') = \sum_{\boldsymbol{h}} \sum_{\boldsymbol{h}'} p(\boldsymbol{h} | \boldsymbol{x}) p(\boldsymbol{h}' | \boldsymbol{x}') K_{z}(\boldsymbol{z}, \boldsymbol{z}')$$
$$= \sum_{\boldsymbol{h}}^{n_{x}} \sum_{\boldsymbol{h}'}^{n_{h}} \gamma_{kl}(\boldsymbol{x}) \gamma_{kl}(\boldsymbol{x}')$$

 $\gamma_{kl}$  is a marginalized count  $\gamma_{kl}(\boldsymbol{x}) = \sum_{\boldsymbol{h}} p(\boldsymbol{h}|\boldsymbol{x}) c_{k\ell}(\boldsymbol{z})$ 

 $k=1 \ell=1$ 

# Computing Marginalized Counts from HMM

Marginalized count is described as

$$\gamma_{kl}(\boldsymbol{x}) = \frac{1}{m} \sum_{i=1}^{m} \sum_{h_i=1}^{n_h} p(h_i | \boldsymbol{x}) I(x_i = k, h_i = \ell).$$

Posterior probability of i-th hidden variable is efficiently computed as

$$p(h_i = k | \boldsymbol{x}) = f_k(i) b_k(i) / p(\boldsymbol{x})$$

 $f_k(i)$ : forward variable,  $b_k(i)$ : backward variable

## 2<sup>nd</sup> order marginalized count kernel

- If adjacent relations between symbols have essential meanings, the count kernel is obviously not sufficient
- 2<sup>nd</sup> order marginalized count kernel
  - 4 neighboring symbols (i.e. 2 visible and 2 hidden) are combined and counted

h: 
$$1 2 2 1 2 2 1 2 2$$
  
x:  $A C G G T T C A A$ 

## **Fisher Kernel**

### Probabilistic mode $p(\boldsymbol{x}|\hat{\boldsymbol{\theta}}), \, \boldsymbol{x} \in \mathcal{X}$

 $\hat{\boldsymbol{\theta}} \in \Re^r$  is a parameter vector obtained from training samples

# Fisher Kernel Mapping to a feature vector (Fisher score vector) $s(x, \hat{\theta}) = \left(\frac{\partial \log p(x|\hat{\theta})}{\partial \theta_1}, \dots, \frac{\partial \log p(x|\hat{\theta})}{\partial \theta_p}\right)$ Inner product of Fisher scores $K_f(x, x') = s(x, \hat{\theta})^\top Z^{-1}(\hat{\theta}) s(x', \hat{\theta})$

Z: Fisher information matrix

# Fisher Kernel from HMM

Derive FK from HMM (Jaakkola et al. 2000) Derivative for emission probabilities only No Fisher information matrix FK from HMM is a special case of marginalized kernels Counts are centralized and weighted  $K_f(\boldsymbol{x}, \boldsymbol{x}') = \sum \sum p(\boldsymbol{h}|\boldsymbol{x}) p(\boldsymbol{h}'|\boldsymbol{x}') K_{fz}(\boldsymbol{z}, \boldsymbol{z}')$ h h' $K_{fz}(z, z') = \sum_{k=1}^{n_x} \sum_{\ell=1}^{n_h} \frac{1}{e_{\ell k}^2} (c_{k\ell}(z) - \bar{c}_{k\ell}(z)) (c_{k\ell}(z') - \bar{c}_{k\ell}(z'))$ 

# Difference between FK and Marginalized Kernels

 FK: Probabilistic model determines the joint kernel and the posterior probabilities

MK: You can determine them separately
 More flexible design!

# Protein clustering experiment

84 proteins containing five classes
 gyrB proteins from five bacteria species
 Clustering methods
 HMM + {FK,MCK1,MCK2}+K-Means
 Evaluation
 Adjusted Rand Index (ARI)

# **Kernel Matrices**











-0.02

0.04

0.02

0

-0.02

30

0.04

MCK2





0.04

0.02

0

-0.02



# **Clustering Evaluation**



# Applications since then..

- Marginalized Graph Kernels (Kashima et al., ICML 2003)
- Sensor networks (Nyugen et al., ICML 2004)
- Labeling of structured data (Kashima et al., ICML 2004)
- Robotics (Shimosaka et al., ICRA 2005)
- Kernels for Promoter Regions (Vert et al., NIPS 2005)
- Web data (Zhao et al., WWW 2006)

## Summary (Marginalized Kernels)

 General Framework for using generative model for defining kernels
 Fisher kernel as a special case
 Broad applications

# Part 3 Marginalized Graph Kernels

## **Motivations for graph analysis**

### Existing methods assume " tables"

 Serial Num	Name	Age	Sex	Address	
 0001	00	40	Male	Tokyo	
0002	××	31	Female	Osaka	-

### Structured data beyond this framework

### → New methods for analysis

AQFERTL		IVNEYS	Y	I	VYLEGCT	P. knowlesi
AQFERTL	L	IVNEYS	Y	I	VYLEGCT	P. simiovale
AQFERTL	L	IVNEYS	Y	I	VYLEGCT	P.v./chesson
AQFERTL	L	IVNEYS	H	I	VYLEGCT	P. simium
AQFERTL	L	IVNEYS	н	I	VYLEGCT	P.v./Africa
AQFERTL	L	IVNEYS	н	I	VYLEGCT	P.v./Thai-1090
AQFERTL	L	IVNEYS	н	I	VYLEGCT	P.v./Thai-115
AQFERTL	L	IVNEYS	н	I	VYLEGCT	P.v./N.Korea
AQFERTL	L	IVNEYS	н	I	VYLEGCT	P.v./Vietnam
AQFERTL	L	IVNEYS	н	v	VYLEGCT	P.v./Salvador-1
AQFERTL	L	IVNEYS	н	v	VYLEGCT	P.v./Salvador-2
AQFERTL	L	IVNEYS	н	v	VYLEGCT	P.v./Brazil-1
AQFERTL	L	IVNEYS	н	v	VYLEGCT	P.v./Brazil-2
AQFERTL	L	IVNEYS	н	v	VYLEGCT	P.v./Honduras-1
AQFERTL	L	IVNEYS	н	v	VYLEGCT	P.v./Honduras-2
AQFERTL	L	IVNEYS	н	v	VYLEGCT	P.v./Panama
		,		-	,	













































### Marginalized Graph Kernels (Kashima, Tsuda, Inokuchi, ICML 2003)

### • Going to define the kernel function K(G, G')• Both vertex and edges are labeled



## Label path

Sequence of vertex and edge labels  $\boldsymbol{h} = (A, e, A, d, D, a, B, c, D)$ Generated by random walking Uniform initial, transition, terminal probabilities b D а С Ô e

B

# Path-probability vector

Label path hProbability p(h|G)AaA0.001::AcDbE0.000003::AeAdDaBcD0.0000007

# Kernel definition

## Kernels for paths

 $K(h, h') = \begin{cases} 0 & (|h| \neq |h'|) \\ k_v(h_1, h'_1)k_e(h_2, h'_2) \cdots k_v(h_\ell, h'_\ell) & (|h| = |h'|) \end{cases}$ 

Take expectation over all possible paths!
 Marginalized kernels for graphs
 K(G,G') = \sum\_h \sum\_{h'} p(h|G)p(h'|G')K(h,h')

A C D b E ↓ ↓ ↓ ↓ B C D a A
#### Computation

#### Transition probability : $\lambda$

Initial and terminal : omitted

•  $S_v(x)$  : Set of paths ending at  $\nu$ 

•  $K_V$ : Kernel computed from the paths ending at (v, v')

$$K_V(v,v') = \sum_{s \in S_v(x)} \sum_{s' \in S_{v'}(x')} \lambda^{|s|} \lambda^{|s'|} K_S(s,s')$$

•  $K_V$  is written recursively

$$K_V(v,v') = \lambda^2 I(v=v') \Big( 1 + \sum_{\tilde{v} \in A(v)} \sum_{\tilde{v}' \in A(v')} \lambda^2 K_V(\tilde{v},\tilde{v}') \Big)$$

 Kernel computed by solving linear equations (polynomial time)



## **Graph Kernel Applications**

Chemical Compounds (Mahe et al., 2005)
 Protein 3D structures (Borgwardt et al, 2005)
 RNA graphs (Karklin et al., 2005)
 Pedestrian detection
 Signal Processing

#### **Predicting Mutagenicity**

MUTAG benchmark dataset
 Mutation of Salmonella typhimurium
 125 positive data (effective for mutations)
 63 negative data (not effective for mutations)

 Table 5: Accuracy Results Obtained for the Leave-One-Out

 Classification of the "Unfriendly Part" of the Mutag Data Set<sup>a</sup>

Lin.Reg	Lin.Reg+	DT	NN	Progol1	Progol2	graph kernels
66.7%	71.8%	83.3%	69.0%	85.7%	83.3%	88.1%

<sup>a</sup> Lin.Reg (Linear Regression), DT (Decision Tree), NN (Neural Network), and Progol1/2 (Inductive Logic Programming): ref 19.

Mahe et al. J. Chem. Inf. Model., 2005

# Classification of Protein 3D structures

#### Graphs for protein 3D structures

- Node: Secondary structure elements
- Edge: Distance of two elements
- Calculate the similarity by graph kernels



Borgwardt et al. "Protein function prediction via graph kernels", ISMB200545

## Classification of proteins: Accuracy

Kernel type	Accuracy	St. dev.
Vector kernel	76.86	1.23
Optimized vector kernel	80.17	1.24
Graph kernel	77.30	1.20
Graph kernel without structure	72.33	5.32
Graph kernel with global info	84.04	3.33
DALI classifier	75.07	4.58

**Table 1.** Accuracy of prediction of functional class of enzymes and nonenzymes in 10-fold cross-validation with C-SVM. The first two results are the results obtained by Dobson and Doig (2003). "Graph kernel" is our protein kernel defined as in Section 2.3, "Graph kernel without structure" is the same kernel but on protein models without structural edges, "Graph kernel with global info" is our protein graph kernel plus additional global node labels. "DALI classifier" is a Nearest Neighbor Classifier on DALI Z-scores.

Borgwardt et al. "Protein function prediction via graph kernels", ISMB200546

#### Pedestrian detection in images (F. Suard et al., 2005)







# Classifying RNA graphs (Y. Karklin et al.,, 2005)



## Strong points of MGK

Polynomial time computation O(n^3)
Positive definite kernel

Support Vector Machines
Kernel PCA
Kernel CCA
And so on...

#### Drawbacks of graph kernels

Global similarity measure
 Fails to capture subtle differences
 Long paths suppressed
 Results not interpretable

Structural features ignored (e.g. loops)
 No labels -> kernel always 1

## Part 4. Weisfeiler Lehman

## kernel

#### Convert a graph into a set of words



i) Make a label set of adjacent vertices ex) {E,A,D} ii) Sort ex) A,D,E iii) Add the vertex label as a prefix ex) B,A,D,E iv) Map the label sequence to a unique value ex) B,A,D,E $\rightarrow$ R

v) Assign the value as the new vertex label







Courtesy K. Borgwardt

5,23

2,3

3,245

G

Dataset	MUTAG	N	CI1	NC	CI109	D	& D
Maximum # nodes 28		1	11	111		5748	
Average # nodes	17.93	29.87		29.68		284.32	
# labels	7		37		54		89
Number of graphs	188	100	4110	100	4127	100	1178
Weisfeiler-Lehman	6"	5"	7'20"	5"	7'21"	58"	11'
Ramon & Gärtner	40'6"	25'9"	29 days*	26'40"	31 days*		
Graphlet count	3"	2"	1'27"	2"	1'27"	2'40"	30'21"
Random walk	12"	58'30"	68 days*	2h 9'41"	153 days*		
Shortest path	2"	3"	4'38"	3"	4'39"	58'45"	23h 17'2"

----: did not finish in 2 days, \* = extrapolated.

Table 2: CPU runtime for kernel computation on graph classification benchmark datasets

## Part 5. Reaction Graph Kernels

#### KEGG lysine degradation pathway

LYSINE DEGRADATION



00310 8/28/09 (c) Kanehisa Laboratories

#### Missing enzymes in metabolic networks

- Many enzymatic reactions whose substrate and product are known, but the enzyme involved is unknown.
- Need to assign Enzymatic Classification numbers.



## EC number

- EC (Enzymatic Classification) number is a hierarchical categorization of
  - Enzymes
  - Enzymatic reactions

class subsubclass EC 1.3.3.-

EC1, Oxidoreductases













EC5, Isomerases





#### Task

Given a pair of substrate and product as a query, find similar reactions in the database



**Similarity measure of reactions is necessary** 

## **Reaction Graph**

- Represent enzymatic reaction as reaction graph
  - Node: Molecules
  - Edge: chemical relation of molecules (main, leave, co-factor, transferase, ligase)
- Reaction graph kernel: Similarity measure of reaction graphs
  - Molecule = Graph of atoms
  - Reaction graph has 'graph of graphs' structure
  - Extension of existing walk-based graph kernel (Kashima et al., 2003)



## Reaction graph kernels (RGK)

- Two-layered kernels on graphs of graphs
  - Node kernel = walk-based graph kernel of molecules
  - Edge kernel = delta kernel of labels
    - "main", "leave", "cofactor", "transferase", "ligase"



## Simplified Settings

- Query might not come in the complete form
- Remove some edges in the database entries



#### RPAIR

Use only reactant edges (main, leave)

**main-only** Use "main" only Automatic classification of enzymatic reactions in KEGG

- KEGG/REACTION database
  - 4610 reactions with known EC number
  - 6 classes, 50 subclasses, 124 subsubclasses
- Construct nearest neighbor classifier based on the reaction graph kernel
  - Three different levels: class, subclass, subclass
  - Three kernel versions: full, RPAIR, main-only
- Measure leave-one-out classification error

### Prediction Accuracy



- As expected, classification is easier for upper categories, but difficult for lower categories such as subsubclass.
- The order of accuracy (full-edge > RPAIR > main-pair) suggests that detailed edge information contributes to further accuracy.

# Predicting unannotated reactions in plant secondary metabolism

- KEGG pathway "Biosynthesis of Secondary Metabolites"
- Out of unannotated 56 reactions, we have manually assigned ECs of 36 reactions under chemists' guidance
- Comparison with an existing rule-based method: ezyme (Kotera et al., J. Am. Chem. Soc, 2004).
- RGK's accuracy was better than e-zyme. (50% improvement for the top candidate)

#### Case 1: EC 3.1.1



Structures of C02046 and C01479 are almost same except structural isomerism. A hydrolysis occur at a carboxylic-esther bond.

method	rank1	rank2	rank3
RGK	3.1.1	1.14.11	1.14.11
e-zyme	6.1.1	3.1.1	NA



After removing a methyl group of C06175 (2.1.1), C01735 will be produced by oxidation of CH-OH group (1.1.1). In the second reactions, enzymes usually use NAD+/NADP+ as an acceptor.

method	rank1	rank2	rank3
RGK	1.1.1	1.1.1	1.1.1
e-zyme	2.1.1	1.13.12	1.14.14

#### Case 3: EC1.14+2.4.1+5.2.1+3.2.1



This reaction is thought to be a set of 5 reactions by analogy with the pathway from  $C00423 \Rightarrow C01772 \Rightarrow C05158 \Rightarrow C05839 \Rightarrow C05838$ , and to C05851. The last reaction is spontaneous and not enzymatic.

method	rank1	rank2	rank3
RGK	1.14.13	1.2.3	1.2.1
e-zyme	NA	NA-	NA-

#### A difficult case



#### A very similar reaction (below) is found by manual inspection



Multi-step reaction is difficult to analyze. We don't know how many steps are hidden between given substrate and product.

#### **Concluding Remarks: New Frontier**

Developing Algorithms for Learning from Graphs Taming Combinatorial Explosion Recursive Fixed Point Iteration: Graph Kernels Statistical Pruning in Search Tree: Graph Boosting Hashing to a set of words: WL Kernel New ideas are necessary to get beyond the current level of speed and prediction accuracy Deeper integration of learning and mining



#### Reference

[1] H. Kashima, K. Tsuda, and A. Inokuchi, "Marginalized kernels between labeled graphs," *ICML*, pp. 321–328, 2003.

[2] H. Saigo, M. Hattori, H. Kashima, and K. Tsuda, "Reaction graph kernels predict EC numbers of unknown enzymatic reactions in plant secondary metabolism.," *BMC bioinformatics*, vol. 11 Suppl 1, no. Suppl 1, p. S31, Jan. 2010.

[3] N. Shervashidze, P. Schweitzer, V. Leeuwen, E. Jan, K. Mehlhorn, and K. Borgwardt, "Weisfeiler-Lehman graph kernels," *Journal of Machine Learning Research*, vol. 12, pp. 2539–2561, 2011.

[4] K. Tsuda, T. Kin, and K. Asai, "Marginalized kernels for biological sequences," *Bioinformatics*, vol. 18, no. Suppl 1, pp. S268–S275, Jul. 2002.