

A tool for evaluating strategies for grouping of biological data

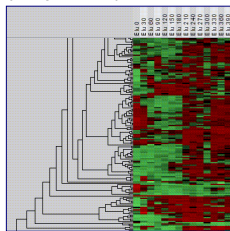
Vaida Jakonienė, Patrick Lambrix

Outline

- Motivation
- Method for similarity based grouping
- KitEGA – illustration
- Summary and future work

Tools for biological data analysis

Hierarchical microarray clustering (J-Express Pro)



Classification of abstracts

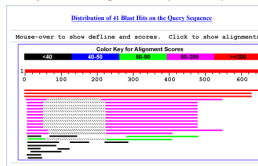
GoPubMed
Ontology-based literature search

levamisole inhibitor [100/135]

- Gene Ontology [37]
- molecular_function [91]
 - catalytic activity [89]
 - signal transducer activity [14]
 - binding [29]
 - antioxidant activity [2]
 - enzyme regulator activity [15]
 - transporter activity [17]
 - cellular_component [25]

Similarity of biological data

Sequence alignment (BLAST)

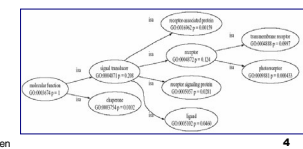


Similarity between data entries

Lord PW, Stevens RD, Brass A, Goble CA. *Bioinformatics*, 19(10):1275-83, 2003.

Molecular Function	Similarity Score
OPSG_HUMAN: Green-sensitive opsin (Green cone photoreceptor pigment)	8.15
OPM_HUMAN: Opn4 (Melanopsin)	7.23
OPSB_HUMAN: Blue-sensitive opsin (Blue cone photoreceptor pigment)	4.92
5H6_HUMAN: 5-hydroxytryptamine 6 receptor (Serotonin receptor)	3.92
A1AA_HUMAN: Alpha-1A adrenergic receptor (Alpha 1A-adrenoceptor)	3.92
A1AB_HUMAN: Alpha-1B adrenergic receptor (Alpha 1B-adrenoceptor)	3.92

Searching with OPSR_HUMAN



- Basic task – computation of a similarity value between objects

Similarity-based grouping

- Not a trivial task
 - data is complex
 - many grouping algorithms available: which algorithm performs best for which grouping task?
 - grouping on which attributes?
 - existing grouping algorithms may not be applied straightforward to new data sets

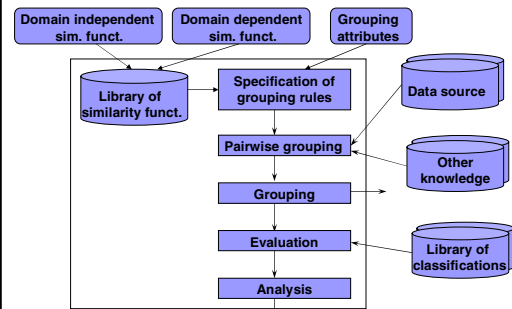
Similarity-based grouping

- Environments that support study, comparison and evaluation of different grouping strategies are needed

Outline

- Motivation
- Method for similarity based grouping
- KitEGA – illustration
- Summary and future work

Method for similarity-based grouping



Outline

- Motivation
- Method for similarity based grouping
- KitEGA – illustration
- Summary and future work

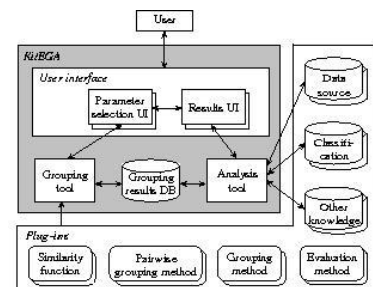


- A toolKit for Evaluating Grouping Algorithms

Idea

- Input components (plug-ins)
 - grouping procedures to be evaluated
 - data sources
 - evaluation methods
 - classifications
 - other knowledge
- Tool executes algorithms and stores results
- User analyzes results using different views on the result data

KitEGA framework



Illustration

- Grouping task. Grouping of proteins with respect to
 - biological function
 - class of isozymes they belong to
- Data source(s)
 - human proteins involved in glycolysis
 - via Entrez retrieved 190 data entries

V. Jakoniene, P. Lambrich. Linköpings universitet, Sweden

13

Data entry

```

LOCUS       NP_000275             390 aa             linear           FRI 16-APR-2006
DEFINITION  pyruvate dehydrogenase (lipoamide) alpha 1 [Homo sapiens].
ACCESSION   NP_000275
VERSION     NP_000275.1   GI:4505695
DBSOURCE    REFSEQ: accession NM_000284.1
KEYWORDS    .
SOURCE      Homo sapiens (human)
  ORGANISM  Homo sapiens
            Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
            Mammalia; Eutheria; Euarchontoglires; Primates; Haplorhini;
            Catarrhini; Hominidae; Homo.
REFERENCE   1 (residues 1 to 390)
  AUTHORS   Hiromasa, Y., Fujisawa, T., Aso, Y. and Roche, T.E.
  TITLE     Organization of the cores of the mammalian pyruvate dehydrogenase
            complex formed by E2 and E2 plus the E3-binding protein and their
            capacities to bind the E1 and E3 components
  JOURNAL   J. Biol. Chem. 279 (9), 6921-6935 (2004)
  PUBMED   14638652
  REMARK    GeneRIF: model of the pyruvate dehydrogenase complex formed by E2
            and E2 plus the E3-binding protein and binding of the E1 and E3
            components
    
```

Entrez. Protein database

V. Jakoniene, P. Lambrich. Linköpings universitet, Sweden

14

Data entry

COMMENT PROVISIONAL REFSEQ: This record has not yet been subject to final NCBI review. The reference sequence was derived from [U33919.1](#).

Summary: The pyruvate dehydrogenase complex is a nuclear-encoded mitochondrial matrix multienzyme complex that provides the primary link between glycolysis and the tricarboxylic acid (TCA) cycle by catalyzing the irreversible conversion of pyruvate into acetyl-CoA. The PDH complex is composed of multiple copies of 3 enzymes: E1 (PDHA1); dihydrolipoyl transacylase (DLAT; MIM 608770) [E1; EC 2.3.1.12]; and dihydrolipoyl dehydrogenase (DLD; MIM 298324) [E2; EC 1.8.1.4]. The E1 enzyme is a heterotetramer of 2 alpha and 2 beta subunits. The E1-alpha subunit contains the E1 active site and plays a key role in the function of the PDH complex (Brown et al., 1994). [supplied by OMIM].

FEATURES

Location/Qualifiers

source

1..390

/organism="Homo sapiens"

/db_xref="taxon:9606"

/chromosome="X"

/map="Xp22.2-p22.1"

Protein

1..390

/product="pyruvate dehydrogenase (lipoamide) alpha 1"

/EC_number="1.3.1.3"

/note="pyruvate dehydrogenase alpha subunit; pyruvate dehydrogenase E1 alpha subunit"

V. Jakoniene, P. Lambrich. Linköpings universitet, Sweden

15

Data entry

```

CDS
1..390
/gene="PDHA1"
/coded_by="NM_000284.1:106..1278"
GO_annot
/go_component="mitochondrion [pmid 3034892]"
/go_function="oxidoreductase activity; oxidoreductase
activity acting on the aldehyde or SMO group of donors,
single as acceptor; pyruvate dehydrogenase
[acetyl-transferase] activity"
/go_process="acetyl-CoA metabolism; glycolysis;
metabolism"
/db_xref="GeneID:5160"
/db_xref="HGNC:8806"
/db_xref="HFPD:0242"
/db_xref="WIKI:300502"
Sequence
ORIGIN
1 mkkmlaavr vlgaaqkpa srivavsrnf endatfeikk cdhrlreepg pvtvtitred
61 gklyrrmqg vrrmelkadq lykqkltgrf chldqgeac ovgleaginp tdhltayra
121 ngftftrgie vrelaeltg rkggoakgkg gsmhmyaknf yggogivgag vplgqiala
181 ckynqkdevo ltlvgdgan gqifeynm aalwklpof lcenthrymz teveraast
241 gyykngdip glrtdgnall vrestafaa ayrcqkpa lmeqtrym gsmadqgva
301 vtrreieqv rskdgpimll kdrrmndaa sveelkaidv evrkteledaa qfatadpep
361 leelgyhiys sdppfevrya ngvikfkvsa
//
    
```

V. Jakoniene, P. Lambrich. Linköpings universitet, Sweden

16

Data sources

DS1: **GO_{ann}** 67 data entries - only terms of GO function ontology analyzed
- only data entries having GO terms

GO Consortium. Mappings between data values and ontological terms:
ec2go – ec_numbers translated into GO terms
spkw2go – swissprot keywords translated into GO terms



V. Jakoniene, P. Lambrich. Linköpings universitet, Sweden

17

Grouping components

- Library of similarity functions
 - EditDist(v_1, v_2)
 - SeqSim(v_1, v_2)
 - SemSim(v_1, v_2) (→ GO ontology)
- Grouping rules
 - GO ontology
- Grouping methods
 - Connected components
 - Cliques

V. Jakoniene, P. Lambrich. Linköpings universitet, Sweden

18

Evaluation methods

- Types of quality measures
 - internal – based on information obtained during the grouping
 - external – with respect to known classes of the grouped data
- In this illustration: external
 - Purity
 - F-measure
 - Entropy
 - Mutual information

V. Jakoniene, P. Lambrich, Linköping universitet, Sweden

19

Classifications

- Manual classification according to
 - biological function
 - classes of isozymes

V. Jakoniene, P. Lambrich, Linköping universitet, Sweden

20

Selection of test case

Data source: Glyc-Funct-AnnEc-onlyGO (DS3)

Grouping rule: SemSim(GOcomb)>0.95

Grouping method: ConnectedComponents

Evaluation method:

- Entropy
- Purity
- MutualInformation
- FMeasure

Source of classes: Glycolysis: by function

next

V. Jakoniene, P. Lambrich, Linköping universitet, Sweden

21

Specification of grouping rules

Data source: Glyc-Funct-AnnEc-onlyGO (DS3)

Grouping rule: SemSim(GOcomb)>0.95

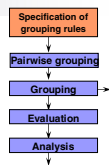
Grouping method: ConnectedComponents

Evaluation method:

- Entropy
- Purity
- MutualInformation
- FMeasure

Source of classes: Glycolysis: by function

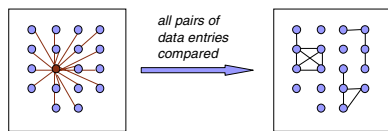
next



V. Jakoniene, P. Lambrich, Linköping universitet, Sweden

22

Pairwise grouping



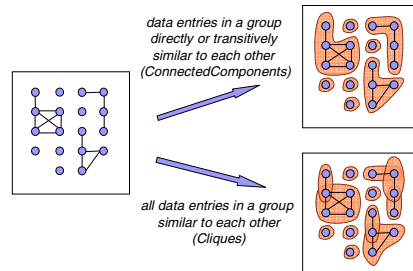
Grouping rule: SemSim(GOcomb)>0.95

Data source: Glyc-Funct-AnnEc-onlyGO (DS3)

V. Jakoniene, P. Lambrich, Linköping universitet, Sweden

23

Grouping



V. Jakoniene, P. Lambrich, Linköping universitet, Sweden

24

Grouping

Specification of grouping rules

Pairwise grouping

Grouping

Evaluation

Analysis

Glyc-Funct-AnnEc-onlyGO + SemSim(GOcomb)>0.95 + ConnectedComponents = Glycolysis by function

GroupNr	ClassNr	ID	Definition	GO combined
0	4	P60174	Triosephosphate isomerase (TIMD) (Triose-phosphate isomerase)	go:0004807
0	4	NP_000356	triosephosphate isomerase 1 [Homo sapiens]	go:0016513, go:0004807
1	2	AAA60068	phosphofructokinase	go:0003872
1	2	NP_001002021	liver phosphofructokinase isoform a [Homo sapiens]	go:0003872
1	2	NP_002617	liver phosphofructokinase isoform b [Homo sapiens]	go:0003872
1	2	NP_002618	phosphofructokinase, platelet [Homo sapiens]	go:0005224, go:0016901, go:000166, go:0016740, go:000287, go:0003872
1	2	NP_000280	phosphofructokinase, muscle [Homo sapiens]	go:0005224, go:0016901, go:000166, go:0016740, go:000287, go:0003872
1	2	P17838	6-phosphofructokinase, liver type (Phosphofructokinase 1) (Phosphofructo-1-kinase)	go:0003872

V. Jakoniene, P. Lambrich. Linköping universitet, Sweden 25

Evaluation

Specification of grouping rules

Pairwise grouping

Grouping

Evaluation

Analysis

Number of entries: 92
 Number of groups: 26
 Number of classes: 25

Entropy: 1.0
 Purity: 1.0
 MutualInformation: 0.8810530832230519
 FMeasure: 0.9939613526570048

- Entropy: average distribution of the data entries in each group among the classes
- Purity: average precision of the groups with respect to their best matching classes
- Mutual information: correspondence on average between each group and class
- F-measure: precision and recall of the classes with respect to their best matching groups on average

V. Jakoniene, P. Lambrich. Linköping universitet, Sweden 26

Analysis

Specification of grouping rules

Pairwise grouping

Grouping

Evaluation

Analysis

Glyc-Funct-AnnEc-onlyGO + SemSim(GOcomb)>0.95 + ConnectedComponents = Glycolysis by function

	0(5)	1(2)	2(14)	3(7)	4(2)	5(4)	6(4)	7(4)	8(4)	9(12)	10(5)	11(4)
0(2)												
1(14)			14/0/0		2/0/0							
2(12)										12/0/0		
3(7)						7/0/0						
4(8)												
5(1)											1/0/4	
6(2)					2/0/0							
7(1)												
8(4)							4/0/0					
9(6)												
10(1)												
11(4)												4/0/1
12(5)	5/0/0											
13(1)												
14(1)												

true positives
 false positives
 false negatives

V. Jakoniene, P. Lambrich. Linköping universitet, Sweden 27

Analysis

group: 11(4) + class: 10(5) + 4 0 1

GroupNr	ClassNr	ID	Definition	GO combined
11	10	P08359	Pyruvate dehydrogenase E1 component alpha subunit, somatic form, mitochondrial precursor (PDHE1-A ty)	go:0004739
11	10	NP_000275	pyruvate dehydrogenase (lipoamide) alpha 1 [Homo sapiens]	go:0016491, go:0004739, go:0016624
11	10	P11177	Pyruvate dehydrogenase E1 component beta subunit, mitochondrial precursor (PDHE1-B)	go:0004739
11	10	P29803	Pyruvate dehydrogenase E1 component alpha subunit, testis-specific form, mitochondrial precursor (PD)	go:0004739
5	10	P10515	[Dihydro]pyruvate acetyltransferase component of pyruvate dehydrogenase complex, mitochondr	go:0004742

V. Jakoniene, P. Lambrich. Linköping universitet, Sweden 28

Analysis - comparison

ID	DataSource	Rule	GrMethod	Classif	# of entries	# of groups	# of classes	Entropy	Purity	MutualInformation	FMeasure
1	Glyc-Funct-AnnEc-onlyGO	SemSim (GOcomb) >0.95	ConnectedComponents	Glycolysis by function	67	26	23	1.0	1.0	0.9117709729626631	0.974650558740109
2	Glyc-Funct-AnnEc-onlyGO	SemSim (GOcomb) >0.95	ConnectedComponents	Glycolysis by function	73	23	24	0.8652854637823463	0.8	0.7942395602417653	0.7917895141899143
3	Glyc-Funct-AnnEc-onlyGO	SemSim (GOcomb) >0.95	ConnectedComponents	Glycolysis by function	92	26	25	1.0	1.0	0.8810230832230519	0.9939613526570048
4	Glyc-Funct-AnnEc-onlyGO	SemSim (GOcomb) >0.85	ConnectedComponents	Glycolysis by function	92	21	25	0.77558996913319	0.6956821799130435	0.6824607500957851	0.70742229746053634
5	Glyc-Funct-AnnEc-onlyGO	SemSim (GOcomb) >0.95	Cliques	Glycolysis by function	92	29	25	1.0	1.0	0.8839495988521302	0.83924244671890822
6	Glyc-Funct-AnnEc-onlyGO	SeqSim(seq) >0.85	ConnectedComponents	Glycolysis by function	92	41	25	1.0	1.0	0.8231661542259041	0.8046256946714613
7	Glyc-Funct-AnnEc-onlyGO	SemSim (GOcomb) >0.95	ConnectedComponents	Glycolysis by function	92	26	47	0.7921820789964245	0.5869565217391305	0.7969682196333567	0.6484704432358992
8	Glyc-Funct-AnnEc-onlyGO	SeqSim(seq) >0.85	ConnectedComponents	Glycolysis by function	92	41	47	0.9256079005200991	0.8478260869166211	0.88481106962638725	0.8380873959690911

- Comparisons: use of different data sources, grouping algorithms, and classifications, grouping on different attributes, impact of threshold

V. Jakoniene, P. Lambrich. Linköping universitet, Sweden 29

Test cases. Observations

- Best suited grouping approaches. For data source Glyc-Funct-AnnEc-onlyGO (DS3)
 - SeqSim (Sequence) for grouping on classes of isozymes
- Suitability of mappings for the used grouping approaches
 - spkw2go – too general, e.g. 'Glycolysis'
 - ec2go – specific enough, e.g. '6-phosphofructokinase activity'

V. Jakoniene, P. Lambrich. Linköping universitet, Sweden 30

Summary and future work

- Motivated need for environments that support the evaluation and comparison of similarity-based grouping procedures
- Implemented the KitEGA tool based on a method for evaluating similarity-based grouping algorithms
- Illustrated KitEGA using test cases based on different strategies and classifications

- Extend the Kitega implementation