

For written notes on this lecture, please read chapter 19 of *The Practical Bioinformatician* and *Hawkins & Kihara, JBCB 5(1):1-30, 2007*

CS2220: Introduction to Computational Biology Lecture 6: Sequence Homology Interpretation

Limsoon Wong



Plan

- Recap of sequence alignment
- Guilt by association
- Active site/domain discovery
- What if no homology of known function is found?
 - Genome phylogenetic profiling
 - SVM-Pairwise
 - Protein-protein interactions
- Key mutation site discovery

Copyright 2010 © Limsoon Wong

Very Brief Recap of Sequence Comparison/Alignment

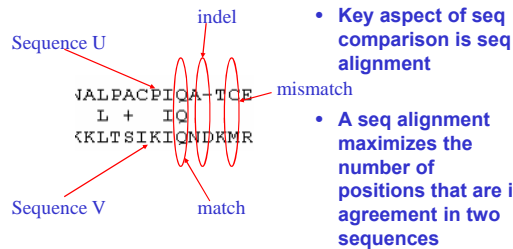


Motivations for Sequence Comparison

- DNA is blue print for living organisms
 - ⇒ Evolution is related to changes in DNA
 - ⇒ By comparing DNA sequences we can infer evolutionary relationships between the sequences w/o knowledge of the evolutionary events themselves
- Foundation for inferring function, active site, and key mutations

Copyright 2010 © Limsoon Wong

Sequence Alignment



Copyright 2010 © Limsoon Wong

Sequence Alignment: Poor Example



- Poor seq alignment shows few matched positions
 - ⇒ The two proteins are not likely to be homologous

Alignment by FASTA of the sequences of amicyanin and domain 1 of ascorbate oxidase

```
Amicyanin      60  70  80  90  100  110  120
WNWNIFFRAGYSGDRAALRQPSRHWKASLFFVQAGEYRFRCTFRERKSHAVY
Ascorbate Oxidase  120  130  140  150  160  170  180
LIRGKIFRAGYSGDRAALRQPSRHWKASLFFVQAGEYRFRCTFRERKSHAVY
```

No obvious match between
Amicyanin and Ascorbate Oxidase

Copyright 2010 © Limsoon Wong

7

Sequence Alignment: Good Example

- Good alignment usually has clusters of extensive matched positions
- ⇒ The two proteins are likely to be homologous

```

>gi|13476732|ref|NP_108301.1| unknown protein [Mesorhizobium loti]
gi|14027493|db|BAB53762.1| unknown protein [Mesorhizobium loti]
Length = 105

Score = 105 bits (262), Expect = 1e-22
Identities = 61/106 (57%), Positives = 73/106 (68%), Gaps = 1/106 (0%)

Query: 1 MKPGRLASIALAIIPLPMVPAHAATIEITMENLVISPTVSAKVGDTIRVWVNDVFAHT 60
          EK G L ++ MA PA AATIE+T+ LV SP V AKVGDIT VWR DV AHT
Sbjct: 1 MEAGALIELSLAALALMAFAAAAATIEVTIDKLVFSPATVEAKVGDITIEVWVNDVFAHT 60
          good match between
          Amicyanin and unknown M. loti protein
  
```

Copyright 2010 © Limsoon Wong

8

Multiple Alignment: An Example

- Multiple seq alignment maximizes number of positions in agreement across several seqs
- seqs belonging to same “family” usually have more conserved positions in a multiple seq alignment

```

gi|126467| FHTSWPFDGVPFPIGLMLKFLKVKACNP--QYAGALVHCSAGVGTGTFVVIDAML
gi|2499753| FHTGWPDHGVPYHATGLLSFIRVKLENP--PSAGLIVHCSAGAGRTGCTYIVIDIMLD
gi|462550| YHQTQWDMGVPYALVPLVTFVRSAAARM--PETGPIVHCSAGVGTGTYIVIDSM
gi|2499751| FHTSWPDHGVPTDILLINFRYLVRDYKOSPPESPIIVHCSAGVGTGTFIAIDRLIY
gi|1709906| FQFTAFPDHGVPEHPTPLAFLERVKTCNP--PDAGPIVHCSAGVGTGCTYIVDAHLE
gi|1264711| LHTSWPFDGVPFPIGLMLKFLKVKLENP--VHAGPIVHCSAGVGTGTYIVIDAHMA
gi|548626| FHTGWPDHGVPYHATGLLSFIRVKLENP--PSAGLIVHCSAGAGRTGCTYIVIDIMLD
gi|1315701| FHTGWPDHGVPYHATGLLSFIRVKLENP--VHAGPIVHCSAGVGTGCTYIVIDIMLD
gi|2144715| FHTSWPFDGVPDITDILLINFRYLVRDYKOSPPESPIIVHCSAGVGTGTFIAIDRLIY
          .. * * * * *
          .. * * * * *
  
```

Conserved sites


Copyright 2010 © Limsoon Wong

Application of Sequence Comparison: Guilt-by-Association

10

A protein is a ...

- A protein is a large complex molecule made up of one or more chains of amino acids
- Protein performs a wide variety of activities in the cell



Copyright 2010 © Limsoon Wong

11

Function Assignment to Protein Sequence

```

SPSTNRKYPPLPVDKLEEEINRRMADNKLFREEFNALPACPIQATCEAASKEENKEKNR
YVNLILPYDHSRVHLTPVEGVPSDYLINASFINGYQEKNFIAAQGPKEETVNDFFRWIWE
QNTATIVMIVNLKKEKCKCAQYWPDQSCWTYGNVRSVVDVTVLVDYTVRKFCCIQQGVD
VTNRKPKRLITQFHTSWPFDGVPFPIGLMLKFLKVKACNPQYAGALVHCSAGVGRGT
TFVVIDAMLDMHSEKRVVDYVGFVRSIRAQRQCMVQTDQYVFIYQALLEHYLYGDTELE
VT
  
```

- How do we attempt to assign a function to a new protein sequence?

Copyright 2010 © Limsoon Wong

12

Invariant and Abductive Reasoning

- Function is determined by 3D struct of protein & environment protein is in
- Constraints imposed by 3D struct & environment give rise to “invariant” properties observed in proteins having that function

⇒ **Abductive reasoning**

- If those invariant properties are seen in a protein, then the protein is homolog of this protein

Entailment $A \rightarrow B$

⇒ **“Guilt by association”**

Copyright 2010 © Limsoon Wong

13

Guilt-by-Association

- Compare the target sequence T with sequences S_1, \dots, S_n of known function in a database
- Determine which ones amongst S_1, \dots, S_n are the mostly likely homologs of T
- Then assign to T the same function as these homologs
- Finally, confirm with suitable wet experiments

Copyright 2010 © Limsoon Wong

14

Guilt-by-Association

Compare T with seqs of known function in a db

Poor Sequence Alignment

- Poor seq alignment shows few matched positions
- The two proteins are not likely to be homologous

Alignment by FASTA of the sequences of *amylomannan* and domain 1 of *aspartate oxidase*

No obvious match between *Ameyman* and *Aspartate Oxidase*

Good Sequence Alignment

- Good alignment usually has clusters of extensive matched positions
- The two proteins are likely to be homologous

Alignment by FASTA of the sequences of *amylomannan* and domain 1 of *aspartate oxidase*

Ameyman and Aspartate Oxidase

Assign to T same function as homologs

Confirm with suitable wet experiments

Discard this function as a candidate

Copyright 2010 © Limsoon Wong

15

BLAST: How It Works

Altschul et al., *JMB*, 215:403-410, 1990

- BLAST is one of the most popular tool for doing "guilt-by-association" sequence homology search

find from db seqs with short perfect matches to query seq

find seqs with good flanking alignment

Exercise: Why do we need this step?

Copyright 2010 © Limsoon Wong

16

Homologs obtained by BLAST

Sequences producing significant alignments:	Score (bits)	E Value
ml14192729 emb AA654109.1 AF331081.1 protein tyrosin phosph...	62.1	e-177
ml1264671 emb F184531 PTA_HUMAN Protein-tyrosine phosphatase...	62.1	e-177
ml14506503 ref PDB_028287.1 protein tyrosine phosphatase, r...	62.0	e-176
ml12372824 ref U1291300A protein Tyr phosphatase	62.0	e-176
ml16450369 ref PDB_543030.1 protein tyrosine phosphatase, ...	61.4	e-176
ml1320671 emb CA437447.1 tyrosine phosphatase precursor [Ho...	61.4	e-176
ml1285113 ipr IP01285 protein-tyrosine-phosphatase (EC 3.1.1...	61.3	e-176
ml16981461 ref PDB_036895.1 protein tyrosine phosphatase, r...	61.1	e-176
ml12029141 emb U17501A Chain A, Receptor Protein Tyrosine Ph...	61.0	e-174
ml132131 emb CA32666.1 protein-tyrosine phosphatase (Homo...	61.0	e-174
ml14506503 emb AA654109.1 protein tyrosine phosphatase >g 4...	60.5	e-172
ml16679587 ref PDB_033006.1 protein tyrosine phosphatase, r...	60.1	e-172
ml14832821 emb AA117990.1 protein tyrosine phosphatase alpha	59.9	e-170

- Thus our example sequence could be a protein tyrosine phosphatase α (PTP α)

Copyright 2010 © Limsoon Wong

17

Example Alignment with PTP α

```

Score = 652 bits (1629), Expect = e-180
Identities = 294/302 (97%), Positives = 294/302 (97%)
Query: 1  SPSTNRKPPPLVWVLEEE INSRMADNKLPRFFNALPACP IQVTEAASGGGGGGGG 60
          SPSTNRKPPPLVWVLEEE INSRMADNKLPRFFNALPACP IQATCEAAS  R
Sbjct: 202 SPSTNRKPPPLVWVLEEE INSRMADNKLPRFFNALPACP IQATCEAAS  261

Query: 61  YVNLFPYDNRVHLTPVEQVPSDVTINASFINGVQKQKFLAAGPKERTVDFPKRINE 120
          YVNLFPYDNRVHLTPVEQVPSDVTINASFINGVQKQKFLAAGPKERTVDFPKRINE 120
Sbjct: 262 YVNLFPYDNRVHLTPVEQVPSDVTINASFINGVQKQKFLAAGPKERTVDFPKRINE 321

Query: 121  QNTATITVYTLNLSKREKCAQVDFDQGVYGVVYVYVPLVDITVYKPCIQVWD 180
          QNTATITVYTLNLSKREKCAQVDFDQGVYGVVYVYVPLVDITVYKPCIQVWD 180
Sbjct: 322 QNTATITVYTLNLSKREKCAQVDFDQGVYGVVYVYVPLVDITVYKPCIQVWD 381

Query: 161  VYNGRQSLITQPIFTSDFDQVDFPTTIGMLKFLKQKACNQYAGAIYVVICSAQWRTD 240
          VYNGRQSLITQPIFTSDFDQVDFPTTIGMLKFLKQKACNQYAGAIYVVICSAQWRTD 240
Sbjct: 382 VYNGRQSLITQPIFTSDFDQVDFPTTIGMLKFLKQKACNQYAGAIYVVICSAQWRTD 441

Query: 241  TPVYIDAMLDNDCIERKVDVYVYVSRIDAGQDQVQVDMQVYVYVYVYVYVYVYVYV 300
          TPVYIDAMLDNDCIERKVDVYVYVSRIDAGQDQVQVDMQVYVYVYVYVYVYVYVYV 300
Sbjct: 442 TPVYIDAMLDNDCIERKVDVYVYVSRIDAGQDQVQVDMQVYVYVYVYVYVYVYVYV 501
  
```

Copyright 2010 © Limsoon Wong

18

Guilt-by-Association: Caveats

- Ensure that the effect of database size has been accounted for
- Ensure that the function of the homology is not derived via invalid "transitive assignment"
- Ensure that the target sequence has all the key features associated with the function, e.g., active site and/or domain

Copyright 2010 © Limsoon Wong

19

Law of Large Numbers

- Suppose you are in a room with 365 other people
- Q: What is the prob that there is a person in the room having the same birthday as you?
- A: $1 - (364/365)^{365} = 63\%$
- Q: What is the prob that a specific person in the room has the same birthday as you?
- A: $1/365 = 0.3\%$
- Q: What is the prob that there are two persons in the room having the same birthday?
- A: 100%

Copyright 2010 © Limsoon Wong

20

Interpretation of P-value

- Seq. comparison progs, e.g. BLAST, often associate a P-value to each hit
- P-value is interpreted as prob that a random seq has an equally good alignment
- Suppose the P-value of an alignment is 10^{-6}
- If database has 10^7 seqs, then you expect $10^7 * 10^{-6} = 10$ seqs in it that give an equally good alignment
- ⇒ Need to correct for database size if your seq comparison prog does not do that!

Note: $P = 1 - e^{-E}$


Exercise: Name a commonly used method for correcting p-value for a situation like this

Copyright 2010 © Limsoon Wong

22

Lightning Does Strike Twice!

- Roy Sullivan, a former park ranger from Virginia, was struck by lightning 7 times
 - 1942 (lost big-toe nail)
 - 1969 (lost eyebrows)
 - 1970 (left shoulder seared)
 - 1972 (hair set on fire)
 - 1973 (hair set on fire & legs seared)
 - 1976 (ankle injured)
 - 1977 (chest & stomach burned)
- September 1983, he committed suicide



Cartoon: Ron Hipschman
Data: David Hand

Copyright 2010 © Limsoon Wong

23

Effect of Seq Compositional Bias

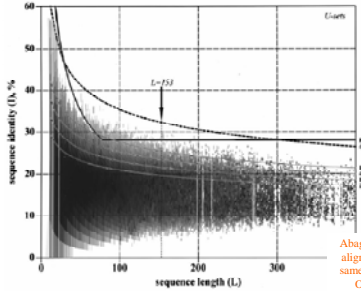
- One fourth of all residues in protein seqs occur in regions with biased amino acid composition
- Alignments of two such regions achieves high score purely due to segment composition
- ⇒ While it is worth noting that two proteins contain similar low complexity regions, they are best excluded when constructing alignments
- E.g., by default, BLAST employs the SEG algo to filter low complexity regions from proteins before executing a search

Source: NCBI

Copyright 2010 © Limsoon Wong

24

Effect of Sequence Length



Abguyan RA, Batalov S. Do aligned sequences share the same fold? J Mol Biol, 1997 Oct 17;273(1):355-68

Copyright 2010 © Limsoon Wong

25

Examples of Invalid Function Assignment: The IMP Dehydrogenases (IMPDH)

15 entries were found

ID	Organism	PIR	Swiss Prot/TrEMBL	RefSeq/GenPept
XP_001181212	Methanococcus jannaschii	MJ0461 unassigned hypothetical protein MJ0461	Y01_MJ216a Hypothetical protein MJ0461	gi 302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii
XP_001181782	Archaeoglobus fulgidus	AF022233 unassigned protein AF022233	Q98441 IMP-DEHYDROGENASE RELATED PROTEIN X	gi 302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii
XP_001181823	Archaeoglobus fulgidus	AF022234 unassigned protein AF022234	Q98442 IMP-DEHYDROGENASE RELATED PROTEIN Y	gi 302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii
XP_001181827	Archaeoglobus fulgidus	AF022235 unassigned protein AF022235	Q98443 IMP-DEHYDROGENASE RELATED PROTEIN Z	gi 302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii
XP_001181828	Archaeoglobus fulgidus	AF022236 unassigned protein AF022236	Q98444 IMP-DEHYDROGENASE RELATED PROTEIN A	gi 302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii
XP_001181829	Archaeoglobus fulgidus	AF022237 unassigned protein AF022237	Q98445 IMP-DEHYDROGENASE RELATED PROTEIN B	gi 302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii
XP_001181830	Archaeoglobus fulgidus	AF022238 unassigned protein AF022238	Q98446 IMP-DEHYDROGENASE RELATED PROTEIN C	gi 302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii
XP_001181831	Archaeoglobus fulgidus	AF022239 unassigned protein AF022239	Q98447 IMP-DEHYDROGENASE RELATED PROTEIN D	gi 302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii
XP_001181832	Archaeoglobus fulgidus	AF022240 unassigned protein AF022240	Q98448 IMP-DEHYDROGENASE RELATED PROTEIN E	gi 302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii
XP_001181833	Archaeoglobus fulgidus	AF022241 unassigned protein AF022241	Q98449 IMP-DEHYDROGENASE RELATED PROTEIN F	gi 302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii
XP_001181834	Archaeoglobus fulgidus	AF022242 unassigned protein AF022242	Q98450 IMP-DEHYDROGENASE RELATED PROTEIN G	gi 302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii
XP_001181835	Archaeoglobus fulgidus	AF022243 unassigned protein AF022243	Q98451 IMP-DEHYDROGENASE RELATED PROTEIN H	gi 302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii
XP_001181836	Archaeoglobus fulgidus	AF022244 unassigned protein AF022244	Q98452 IMP-DEHYDROGENASE RELATED PROTEIN I	gi 302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii
XP_001181837	Archaeoglobus fulgidus	AF022245 unassigned protein AF022245	Q98453 IMP-DEHYDROGENASE RELATED PROTEIN J	gi 302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii
XP_001181838	Archaeoglobus fulgidus	AF022246 unassigned protein AF022246	Q98454 IMP-DEHYDROGENASE RELATED PROTEIN K	gi 302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii
XP_001181839	Archaeoglobus fulgidus	AF022247 unassigned protein AF022247	Q98455 IMP-DEHYDROGENASE RELATED PROTEIN L	gi 302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii g302201667 envs_methanococcus_jannaschii

A partial list of IMP dehydrogenase misnomers in complete genomes remaining in some public databases

Copyright 2010 © Limsoon Wong

26

IMPDH Domain Structure

- Typical IMPDHs have 2 IMPDH domains that form the catalytic core and 2 CBS domains.
- A less common but functional IMPDH (E70218) lacks the CBS domains.
- Misnomers show similarity to the CBS domains

Copyright 2010 © Limsoon Wong

27

Invalid Transitive Assignment

Root of invalid transitive assignment

Mis-assignment of function

No IMPDH domain

Copyright 2010 © Limsoon Wong

28

Emerging Pattern

Typical IMPDH Functional IMPDH w/o CBS

- Most IMPDHs have 2 IMPDH and 2 CBS domains
- Some IMPDH (E70218) lacks CBS domains
- ⇒ IMPDH domain is the emerging pattern

Copyright 2010 © Limsoon Wong

Application of Sequence Comparison: Active Site/Domain Discovery

30

Discover Active Site and/or Domain

- How to discover the active site and/or domain of a function in the first place?
 - Multiple alignment of homologous seqs
 - Determine conserved positions
 - ⇒ Emerging patterns relative to background
 - ⇒ Candidate active sites and/or domains
- Easier if sequences of distance homologs are used

Exercise: Why?

Copyright 2010 © Limsoon Wong


31

In the course of evolution...

Copyright 2010 © Limsoon Wong

32

Multiple Alignment of PTPs



```


g1|126467| FHTSMPDFGVVPT IGLKFLKRVKACNP--QYAGIVVHCSAGVORTGTFVVIDAMLD
g1|2459753| FHTGWFDRGVVPHATGLLSFIRRVKLSNP--PSAGP IUVHCSAGAGRTGCVIVIDIMLD
g1|4625501| YHTTQWDRGVVPEYALPGLTVKRSAAKRP--PETQPLVHCSAGVORTGTFVVIDAMLD
g1|2459751| FHTSMPDRGVVPTDILLINFRVLVDRYKQSPPEPFLVHCSAGVORTGTFVVIDAMLD
g1|1709906| FQTAMPDRGVVPEHTPFLAFLKRVKTCNP--PDAGP IUVHCSAGVORTGTFVVIDAMLD
g1|1264711| LHTSMPDFGVVPT IGLKFLKRVKTLNP--VHAGP IUVHCSAGVORTGTFVVIDAMLD
g1|5486261| FHTGWFDRGVVPHATGLLSFIRRVKLSNP--PSAGP IUVHCSAGAGRTGCVIVIDIMLD
g1|1315701| FHTGWFDRGVVPHATGLLSFRVQVRSKSP--PNAGP IUVHCSAGAGRTGCVIVIDIMLD
g1|2144715| FHTSMPDRGVVPTDILLINFRVLVDRYKQSPPEPFLVHCSAGVORTGTFVVIDAMLD
..* * * * * ..* ..* * * * * ..* * * * * ..* * * * *

```

- Notice the PTPs agree with each other on some positions more than other positions
- These positions are more imp't wrt PTPs
- Else they wouldn't be conserved by evolution
- ⇒ They are candidate active sites


Copyright 2010 © Limsoon Wong

Guilt-by-Association:
What if no homolog of known function is found?



34

What if there is no useful seq homolog?




- **Guilt by other types of association!**
 - Domain modeling (e.g., HMMPFAM)
 - ✓ Similarity of phylogenetic profiles
 - ✓ Similarity of dissimilarities (e.g., SVM-PAIRWISE)
 - Similarity of subcellular co-localization & other physico-chemico properties (e.g., PROTFUN)
 - Similarity of gene expression profiles
 - ✓ Similarity of protein-protein interaction partners
 - ...
 - Fusion of multiple types of info

Copyright 2010 © Limsoon Wong

35

Phylogenetic Profiling

Pellegrini et al., *PNAS*, 96:4285-4288, 1999




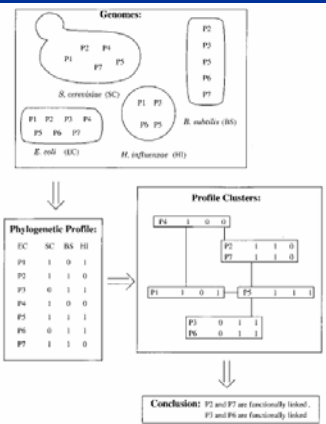
- Gene (and hence proteins) with identical patterns of occurrence across phyla tend to function together

⇒ Even if no homolog with known function is available, it is still possible to infer function of a protein

Copyright 2010 © Limsoon Wong

36


Phylogenetic Profiling: How it Works

Copyright 2010 © Limsoon Wong

37

Phylogenetic Profiling: P-value



The probability of observing by chance z co-occurrences of genes X and Y in a set of N lineages, given that X occurs in x lineages and Y in y lineages is

$$P(z|N, x, y) = \frac{w_z \cdot w_z}{W}$$

where

- $w_z = \binom{N}{z}$
- $w_z = \binom{N-x}{z-x} \cdot \binom{N-y}{z-x}$
- $W = \binom{N}{x} \cdot \binom{N}{y}$

Annotations:

- w_z : No. of ways to distribute z co-occurrences over N lineages
- w_z : No. of ways to distribute the remaining $x-z$ and $y-z$ occurrences over the remaining $N-z$ lineages
- W : No. of ways of distributing X and Y over N lineages without restriction

Copyright 2010 © Limsoon Wong

38

Phylogenetic Profiles: Evidence

Pellegrini et al., *PNAS*, 96:4285-4288, 1999

Keyword	No. of non-homologous proteins in group	No. neighbors in keyword group	No. neighbors in random group
Ribosome	60	197	27
Transcription	36	17	10
tRNA synthase and ligase	26	11	5
Membrane proteins [†]	25	89	5
Flagellar	21	89	3
Iron, ferric, and ferritin	19	21	2
Galactose metabolism	18	21	2
Molybdenin and Molybdenum, and molybdenin	12	6	1
Hypothetical [†]	1,084	198,226	8,440

- E. coli proteins grouped based on similar keywords in SWISS-PROT have similar phylogenetic profiles

Copyright 2010 © Limsoon Wong

39

Phylogenetic Profiling: Evidence

Wu et al., *Bioinformatics*, 19:1524-1530, 2003

fraction of gene pairs having hamming distance D and share a common pathway in KEGG/COG

hamming distance D

hamming distance x,y
 $=$ #lineages X occurs +
 $=$ #lineages Y occurs -
 $2 *$ #lineages X, Y occur

#KEGG
 □ COG

- Proteins having low hamming distance (thus highly similar phylogenetic profiles) tend to share common pathways

Exercise: Why do proteins having high hamming distance also have this behaviour?

Copyright 2010 © Limsoon Wong

40

Guilt by Association of Dissimilarities

Differences of "unknown" to other fruits are same as "apple" to other fruits

	Orange ₁	Banana ₁	...
Apple ₁	Color = red vs orange Skin = smooth vs rough Size = small vs small Shape = round vs round	Color = red vs yellow Skin = smooth vs smooth Size = small vs small Shape = round vs oblong	...
Orange ₂	Color = orange vs orange Skin = rough vs rough Size = small vs small Shape = round vs round	Color = orange vs yellow Skin = rough vs smooth Size = small vs small Shape = round vs oblong	...
Unknown ₁	Color = red vs orange Skin = smooth vs rough Size = small vs small Shape = round vs round	Color = red vs yellow Skin = smooth vs smooth Size = small vs small Shape = round vs oblong	...

"unknown" is an "apple"!

Copyright 2010 © Limsoon Wong

41

SVM-Pairwise Framework

Training Data
 S1
 S2
 S3
 ...

Feature Generation

Training Features
 S_1, S_2, S_3, \dots
 $f_{11}, f_{12}, f_{13}, \dots$
 $f_{21}, f_{22}, f_{23}, \dots$
 $f_{31}, f_{32}, f_{33}, \dots$
 ...

Training

Support Vectors Machine (Radial Basis Function Kernel)

Trained SVM Model (Feature Weights)

Testing Data
 T1
 T2
 T3
 ...

Feature Generation

Testing Features
 T_1, T_2, T_3, \dots
 $f_{11}, f_{12}, f_{13}, \dots$
 $f_{21}, f_{22}, f_{23}, \dots$
 $f_{31}, f_{32}, f_{33}, \dots$
 ...

Classification

RBF Kernel

Discriminant Scores

f_{1j} is the local alignment score between S_j and S_i

f_{1j} is the local alignment score between T_j and S_i

Image credit: Kenny Chua

Copyright 2010 © Limsoon Wong

42

Performance of SVM-Pairwise

- Receiver Operating Characteristic (ROC)
 - The area under the curve derived from plotting true positives as a function of false positives for various thresholds.
- Rate of median False Positives (RFP)
 - The fraction of negative test examples with a score better or equals to the median of the scores of positive test examples.

Copyright 2010 © Limsoon Wong

Protein Function Prediction from Protein Interactions

Level-1 neighbour

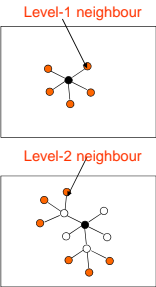
Level-2 neighbour

NUS
 National University of Singapore

44

Functional Association Thru Interactions

- **Direct functional association:**
 - Interaction partners of a protein are likely to share functions w/ it
 - Proteins from the same pathways are likely to interact
- **Indirect functional association**
 - Proteins that share interaction partners with a protein may also likely to share functions w/ it
 - Proteins that have common biochemical, physical properties and/or subcellular localization are likely to bind to the same proteins

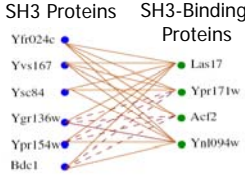


Copyright 2010 © Limsoon Wong

45

An illustrative Case of Indirect Functional Association?

SH3 Proteins SH3-Binding Proteins

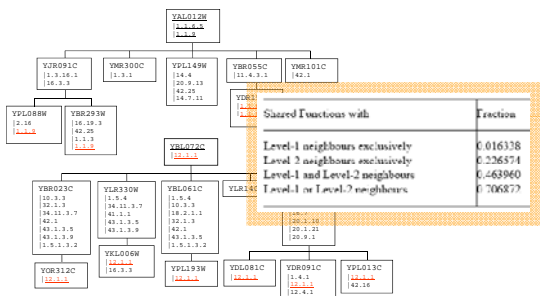


- Is *indirect functional association* plausible?
- Is it found often in real interaction data?
- Can it be used to improve protein function prediction from protein interaction data?

Copyright 2010 © Limsoon Wong

46

Freq of Indirect Functional Association



Shared Functions with	Function
Level-1 neighbours exclusively	0.016318
Level-2 neighbours exclusively	0.226574
Level-1 and Level-2 neighbours	0.463960
Level-1 or Level-2 neighbours	0.706877

Source: Kenny Chua

Copyright 2010 © Limsoon Wong

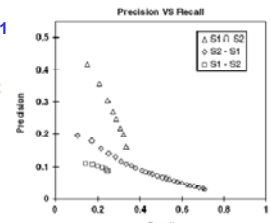
47

Prediction Power By Majority Voting

- Remove overlaps in level-1 and level-2 neighbours to study predictive power of "level-1 only" and "level-2 only" neighbours
- Sensitivity vs Precision analysis

$$PR = \frac{\sum_i k_i}{\sum_i m_i} \quad SN = \frac{\sum_i k_i}{\sum_i n_i}$$

- n_i is no. of fn of protein i
- m_i is no. of fn predicted for protein i
- k_i is no. of fn predicted correctly for protein i



- ⇒ "level-2 only" neighbours performs better
- ⇒ L1 ∩ L2 neighbours has greatest prediction power

Copyright 2010 © Limsoon Wong

48

Functional Similarity Estimate: Czekanowski-Dice Distance

- **Functional distance between two proteins** (Brun et al, 2003)

$$D(u, v) = \frac{|N_u \Delta N_v|}{|N_u \cup N_v| + |N_u \cap N_v|}$$

- N_k is the set of interacting partners of k
- $X \Delta Y$ is symmetric diff betw two sets X and Y
- Greater weight given to similarity

⇒ Similarity can be defined as

$$S(u, v) = 1 - D(u, v) = \frac{2X}{2X + (Y + Z)}$$

Is this a good measure if u and v have very diff number of neighbours?

Copyright 2010 © Limsoon Wong

49

Functional Similarity Estimate: FS-Weighted Measure

- **FS-weighted measure**

$$S(u, v) = \frac{2|N_u \cap N_v|}{|N_u - N_v| + 2|N_u \cap N_v|} \times \frac{2|N_u \cap N_v|}{|N_v - N_u| + 2|N_u \cap N_v|}$$

- N_k is the set of interacting partners of k
- Greater weight given to similarity

⇒ Rewriting this as

$$S(u, v) = \frac{2X}{2X + Y} \times \frac{2X}{2X + Z}$$

Copyright 2010 © Limsoon Wong

50

Correlation w/ Functional Similarity

- Correlation betw functional similarity & estimates

Neighbours	CD-Distance	FS-Weight
S ₁	0.471810	0.498745
S ₂	0.224705	0.298843
S ₁ ∪ S ₂	0.224581	0.29629

- Equiv measure slightly better in correlation w/ similarity for L1 & L2 neighbours

Copyright 2010 © Limsoon Wong

51

Reliability of Expt Sources

- Diff Expt Sources have diff reliabilities
 - Assign reliability to an interaction based on its expt sources (Nabiev et al. 2004)
- Reliability betw u and v computed by:

$$r_{u,v} = 1 - \prod_{i \in E_{u,v}} (1 - r_i)$$
 - r_i is reliability of expt source i
 - E_{u,v} is the set of expt sources in which interaction betw u and v is observed

Source	Reliability
Affinity Chromatography	0.823077
Affinity Precipitation	0.455904
Biochemical Assay	0.666667
Dosage Lethality	0.5
Purified Complex	0.891473
Reconstituted Complex	0.5
Synthetic Lethality	0.37386
Synthetic Rescue	1
Two Hybrid	0.265407

Copyright 2010 © Limsoon Wong

52

Functional Similarity Estimate: FS-Weighted Measure with Reliability

- Take reliability into consideration when computing FS-weighted measure:

$$S_k(u,v) = \frac{2 \sum_{w \in (N_u \cap N_v)} r_{u,w} r_{v,w}}{\left(\sum_{w \in N_u} r_{u,w} + \sum_{w \in N_v} (1 - r_{u,w}) \right) + 2 \sum_{w \in (N_u \cap N_v)} r_{u,w}} \times \frac{2 \sum_{w \in (N_u \cap N_v)} r_{u,w} r_{v,w}}{\left(\sum_{w \in N_u} r_{u,w} + \sum_{w \in N_v} (1 - r_{u,w}) \right) + 2 \sum_{w \in (N_u \cap N_v)} r_{u,w}}$$

- N_k is the set of interacting partners of k
- r_{u,w} is reliability weight of interaction betw u and v

⇒ Rewriting

$$S(u,v) = \frac{2X}{2X+Y} \times \frac{2X}{2X+Z}$$

Copyright 2010 © Limsoon Wong

53

Integrating Reliability

- Equiv measure shows improved correlation w/ functional similarity when reliability of interactions is considered:

Neighbours	CD-Distance	FS-Weight	FS-Weight R
S ₁	0.471810	0.498745	0.532596
S ₂	0.224705	0.298843	0.375317
S ₁ ∪ S ₂	0.224581	0.29629	0.363025

Copyright 2010 © Limsoon Wong

54

Improvement to Prediction Power by Majority Voting

Copyright 2010 © Limsoon Wong

55

Improvement to Over-Rep of Functions in Neighbours

Copyright 2010 © Limsoon Wong

56

Use L1 & L2 Neighbours for Prediction

- FS-weighted Average

$$f_x(u) = \frac{1}{Z} \left[\lambda r_{int} \pi_x + \sum_{v \in N_u} S_{TR}(u, v) \delta(v, x) + \sum_{w \in N_u} S_{TR}(u, w) \delta(w, x) \right]$$

- r_{int} is fraction of all interaction pairs sharing function
- λ is weight of contribution of background freq
- $\delta(k, x) = 1$ if k has function x , 0 otherwise
- N_k is the set of interacting partners of k
- π_x is freq of function x in the dataset
- Z is sum of all weights

$$Z = 1 + \sum_{v \in N_u} S_{TR}(u, v) + \sum_{w \in N_u} S_{TR}(u, w)$$

Copyright 2010 © Limsoon Wong

57

Performance of FS-Weighted Averaging

- LOOCV comparison with Neighbour Counting, Chi-Square, PRODISTIN

Copyright 2010 © Limsoon Wong

58

About the Inventor: Chua Hon Nian

- Chua Hon Nian
 - PhD, NUS, 2008
 - Currently postdoc at Harvard
 - 49th hottest paper in Computer Science published in 2006
 - Winner, DREAM2 challenge PPI subnetwork, 2007



Copyright 2010 © Limsoon Wong

Application of
Sequence Comparison:
Key Mutation Site Discovery



60

Identifying Key Mutation Sites

K.L.Lim et al., JBC, 273:28986–28993, 1998

Sequence from a typical PTP domain D2

```
>gi|00000|PTP-D2
EEEEKELTSIKIIONDEKERTCHLPANKEKKEWFLQIIPYEFNRVLIIPYKRGZENTDFWNASF
IDQVYKQDSYIASQGLLETIEDFUREIFRZUSKCSIVHLTELEERQKRCAGYFSPDGLV
SYQDITVLEKKEEKESYTVRDLVFNTRZKESDQIRQVHFEGUPVFGIPIPSDGRHISII
LAVQKQQCSGHPITVCSAGLQRTQYCALSTVLEWVRLDGLDFFQTVKSLRLQRPB
HWQTLQYEPCKYKVVQYVIDAFSDYANPK
```

- Some PTPs have 2 PTP domains
- PTP domain D1 is has much more activity than PTP domain D2
- Why? And how do you figure that out?

Copyright 2010 © Limsoon Wong

61

Emerging Patterns of PTP D1 vs D2

- Collect example PTP D1 sequences
- Collect example PTP D2 sequences
- Make multiple alignment A1 of PTP D1
- Make multiple alignment A2 of PTP D2
- Are there positions conserved in A1 that are violated in A2?
- These are candidate mutations that cause PTP activity to weaken
- Confirm by wet experiments

Copyright 2010 © Limsoon Wong

62

Emerging Patterns of PTP D1 vs D2

This site is consistently conserved in D1, but is not consistently missing in D2
 ⇒ it is not an EP
 ⇒ not a likely cause of D2's loss of function

Exercise: Why?

This site is consistently conserved in D1, but is consistently missing in D2
 ⇒ it is an EP
 ⇒ possible cause of D2's loss of function

absent (red X)
 present (red oval)

Copyright 2010 © Limsoon Wong

63

Key Mutation Site: PTP D1 vs D2

```

g1|00000|P D2 QFHFHGVPEVGIIPSDGKGNISIIAAVQKQQQ--SNGHP IVVHCSAGAGRTGTFCAI?STVL?
g1|1264671 QFHTSVPDFGVVPTPIGLKFLKVKACNP--QIAGIVVHCSAGVGRGTGTFVIVDAML
g1|2499753 QFHTGVPDHGVPYHATGLLSP IRRVLSNP--PSAGP IVVHCSAGAGRTGCVIVIDIML
g1|4625501 QYHVTQVPDHGVPYALFVLTFRSSAAARM--PETGPIVHCSAGVGRGTGTYIVIDSML
g1|2499751 QFHTSVPDHGVPDTELLINFRVLRVYMKQSPFESP ILVHCSAGVGRGTGTF IAI?DRLI
g1|1709906 D1 QFQFTA!VPDHGVP!PH?TFELFLRRVKT?CNP--PDAGPIVHCSAGVGRGTGCF IVIDAML
g1|1264711 QLHFTSVPDFGVVPTPIGLKFLKVK!TLNP--VHAGP IVVHCSAGVGRGTGTF IVIDAMM
g1|5486261 QFHTGVPDHGVPYHATGLLSP IRRVLSNP--PSAGP IVVHCSAGAGRTGCVIVIDIML
g1|1315701 QFHTGVPDHGVPYHATGLLSPVQVSKSP--PHAGPIVHCSAGAGRTGCF IVIDIML
g1|2144715 QFHTSVPDHGVPDTELLINFRVLRVYMKQSPFESP ILVHCSAGVGRGTGTF IAI?DRLI
*.. **.*.*
  
```

- Positions marked by "!" and "?" are likely places responsible for reduced PTP activity
 - All PTP D1 agree on them
 - All PTP D2 disagree on them

Copyright 2010 © Limsoon Wong

64

Key Mutation Site: PTP D1 vs D2

```

g1|00000|P D2 QFHFHGVPEVGIIPSDGKGNISIIAAVQKQQQ--SNGHP IVVHCSAGAGRTGTFCAI?STVL?
g1|1264671 QFHTSVPDFGVVPTPIGLKFLKVKACNP--QIAGIVVHCSAGVGRGTGTFVIVDAML
g1|2499753 QFHTGVPDHGVPYHATGLLSP IRRVLSNP--PSAGP IVVHCSAGAGRTGCVIVIDIML
g1|4625501 QYHVTQVPDHGVPYALFVLTFRSSAAARM--PETGPIVHCSAGVGRGTGTYIVIDSML
g1|2499751 QFHTSVPDHGVPDTELLINFRVLRVYMKQSPFESP ILVHCSAGVGRGTGTF IAI?DRLI
g1|1709906 D1 QFQFTA!VPDHGVP!PH?TFELFLRRVKT?CNP--PDAGPIVHCSAGVGRGTGCF IVIDAML
g1|1264711 QLHFTSVPDFGVVPTPIGLKFLKVK!TLNP--VHAGP IVVHCSAGVGRGTGTF IVIDAMM
g1|5486261 QFHTGVPDHGVPYHATGLLSP IRRVLSNP--PSAGP IVVHCSAGAGRTGCVIVIDIML
g1|1315701 QFHTGVPDHGVPYHATGLLSPVQVSKSP--PHAGPIVHCSAGAGRTGCF IVIDIML
g1|2144715 QFHTSVPDHGVPDTELLINFRVLRVYMKQSPFESP ILVHCSAGVGRGTGTF IAI?DRLI
*.. **.*.*
  
```

Positions marked by "!" are even more likely as 3D modeling predicts they induce large distortion to structure

Copyright 2010 © Limsoon Wong

65

Confirmation by Mutagenesis Expt

- What wet experiments are needed to confirm the prediction?
 - Mutate E → D in D2 and see if there is gain in PTP activity
 - Mutate D → E in D1 and see if there is loss in PTP activity

Exercise: Why do you need this 2-way expt?

Copyright 2010 © Limsoon Wong

66

About the Inventor: Prasanna Kolatkar

- Prasanna Kolatkar**
 - Research Fellow, BIC, NUS, 1997-1999
 - Currently Group Leader at GIS

Copyright 2010 © Limsoon Wong

Concluding Remarks

What have we learned?



- **General methodologies & applications**
 - Guilt by association for protein function inference
 - Invariants for active site discovery
 - Emerging patterns for mutation site discovery
- **Important tactics**
 - Genome phylogenetic profiling
 - SVM-Pairwise
 - Protein-protein interactions

Any Question?



Acknowledgements



- **Some of the slides are based on slides given to me by Kenny Chua**

References



- T.F.Smith & X.Zhang. "The challenges of genome sequence annotation or "The devil is in the details", *Nature Biotech*, 15:1222--1223, 1997
- D. Devos & A.Valencia. "Intrinsic errors in genome annotation", *TIG*, 17:429--431, 2001
- K.L.Lim et al. "Interconversion of kinetic identities of the tandem catalytic domains of receptor-like protein tyrosine phosphatase PTP-alpha by two point mutations is synergist and substrate dependent", *JBC*, 273:28986--28993, 1998
- S.F.Altshul et al. "Basic local alignment search tool", *JMB*, 215:403--410, 1990
- S.F.Altshul et al. "Gapped BLAST and PSI-BLAST: A new generation of protein database search programs", *NAR*, 25(17):3389--3402, 1997

References



- S.E.Brenner. "Errors in genome annotation", *TIG*, 15:132--133, 1999
- M. Pellegrini et al. "Assigning protein functions by comparative genome analysis: Protein phylogenetic profiles", *PNAS*, 96:4285--4288, 1999
- J. Wu et al. "Identification of functional links between genes using phylogenetic profiles", *Bioinformatics*, 19:1524--1530, 2003
- L.J.Jensen et al. "Prediction of human protein function from post-translational modifications and localization features", *JMB*, 319:1257--1265, 2002
- C. Wu, W. Barker. "A Family Classification Approach to Functional Annotation of Proteins", *The Practical Bioinformatician*, Chapter 19, pages 401--416, WSPC, 2004

References



- H.N. Chua, W.-K. Sung. [A better gap penalty for pairwise SVM](#). Proc. APBC05, pages 11-20
- Hon Nian Chua, Wing Kin Sung, Limsoon Wong. [Exploiting Indirect Neighbours and Topological Weight to Predict Protein Function from Protein-Protein Interactions](#). *Bioinformatics*, 22:1623-1630, 2006.
- T. Jaakkola, M. Diekhans, and D. Haussler. A discriminative framework for detecting remote homologies. *JCB*, 7(1-2):95--11, 2000
- T. Hawkins and D. Kihara. Function prediction of uncharacterized proteins. *JBCB*, 5(1):1-30, 2007