

Lecture – 13b

Non-redundant Datasets II

Clustering Methods Based On Composition

Hamming distance

$D^H = \sum |Comp(1)_i - comp(2)_i|, i=1,20$

Comp(i) = $n \times \frac{i}{N}$

A C D E F F G H I K L M N P A R S T V W Y

1. ADIKLAAIKL \rightarrow 03 - 01 - - - - 02 02 - - - - -
2. ADSKLAAIKA \rightarrow 04 - 01 - - - - 01 01 - - - - 01 - -
3. KILASDPQWE \rightarrow 01 - 01 01 - - - 01 01 - - 01 01 - 01 -

$D^H(1-2) = 0.4$
 $D^H(2-3) = 0.8$
 $D^H(1-3) = 1.0$

0.09
 0.01 x 5 = 0.05
 0.05 x 5 = 0.25
 0.25 x 5 = 1.25
 1.25 x 5 = 6.25

M. Michael Gromiha, NPTEL, Bioinformatics, Lecture 13

So, what is the different composition for a sequence 1? How many residues? 1, 2, 3, 4, 5, 6, 7, 8, 9, 10; for simplicity I put 10, but usually what is average length of protein sequences?

Student: 100.

Right you can see large crowd around 200 to 300 or 100 to 300; average length is 315 residues depending upon the sequences, deposited in the UniProt database. Also I showed the graph connecting the number of residues versus the frequency. You must know 100 to 300 residues; so, for example if we take; for simplicity I use small peptides what is the composition for A in the first sequence?

Student: 0.3.

0.3; C, 0; D 0.1; E?

Student: 0.

F?

Student: 0

0; G?

Student: 0

0. H 0; I?

Student: 0.2.

0.2; K?

Student: 0.2.

K is 0.2.

Student: L is also 0.2

L? 0.2; M no; N?

Student: 0.

No.

Student: All other 0.

All are 0; 0.2, 0.4, 0.6, 0.91; then for the second one?

Student: 0.4

0.4.

Student: 0.

0.1.

Student: 0

E, F, G, H, I is 0.1; K?

Student: 0.2

0.2; L?

Student: 0.1.

0.1; 4, 5, 6, 8, 9 there is one more, S; 0.1. So, third sequence; so, what is the composition for alanine?

Student: 0.1

0.1; no C; D 0.1; E 0.1; F no, G; H, I 0.1; K, 0.1; L 0.1 M, N, P 0.1; Q 0.1; no R, no S.

Student: 0.1.

W 0.1; 1, 2, 3, 4, 5, 6, 7, 8, 9; one is missing S 0.1. So, now we calculate Hamming distance between 1 and 2; what is the value Hamming distance between 1 and 2?

Student: 0.1? 0.2, 0.3.

Point?

Student: 0.3.

0.3.

Student: Hm.

This is 0.1, 0, 0.2, 0.3, 0.4; now D, H between 1 and 2 and 3; 0.8.

Student: 0.8.

Very high; so between the 1 and 3; 1.0; so, we have the different sequences; something is wrong right. So, three different sequences; so we can calculate the Hamming distance between any pairs. Now, if you see this one; which two sequences are close to each other?

Student: 1 and 2.

1 and 2 are close to each other; if you see this one; 1 and 2, they are close to each other; we get the value as 0.4. So, we can see the different sequences; then if you have any specific threshold for example, if we have threshold of 0.5; then 1 and 2 are close to each other; either 1 or 2. Then we can take 3; you can make into different clusters and from each cluster, we can take a representative data right and finally, you get the non redundant set of sequences; this one is Hamming distance.

(Refer Slide Time: 05:18)

Clustering Methods Based On Composition

Euclidean distance

$$D^E = \{\sum [\text{Comp}(1)_i - \text{comp}(2)_i]^2\}^{1/2}$$

1. ADIKLAAIKL
2. ADSKLAAIKA
3. KILASDPQWE

$D(1,2) = \sqrt{0.01 + 0.01 + 0.01 + 0.01} = \sqrt{0.04}$

$D(1,2) = 0.2$

$D(2,3) = \sqrt{0.14} = 0.37$

$D(1,3) = \sqrt{0.11} = 0.33$

M. Michael Gromiha, NPTEL, Bioinformatics, Lecture 13

So, another is similar to Hamming distance; we can also get the Euclidean distance. The difference between Hamming distance and Euclidean distance is; what is the difference?

Student: Sum of squares.

Is sum of squares, so you can take composition of 1, minus composition 2; the whole squared, because we do not care about the sign. Then you have to take the root; so what is the value between 1 and 2?

Student: 0.4

That you get the difference right; 0.1, 0.2, 0.3, 0.4; then you have to take the root.

Student: 0.1.

0.1 squares; that is equal to 0.04; 0.1 square equal to 0.01.

Student: Yes.

Plus 0.01 plus 0.01 plus 0.01 and take the root; this is equal to 0.04.

Student: Point.

0.2; this D^E of 1, 2; how about D^E of 2, 3? 2 and 3 is 0.09; 0.01; here it is 0, here it is 0.01; 1, 2 3 times. So; that means, 0.01 into 5; that is equal to 0.09 and 0.05, this is equal to 0.14; so, root of 0.14; this is equal to 0.4? 0.3?

Student: 0.3 something

0.3 something; so, likewise D^E of 1,3; what is 1,3?

Student: 0.4... 0.04

0.04; is 0; this is.

Student: 0.01.

1.

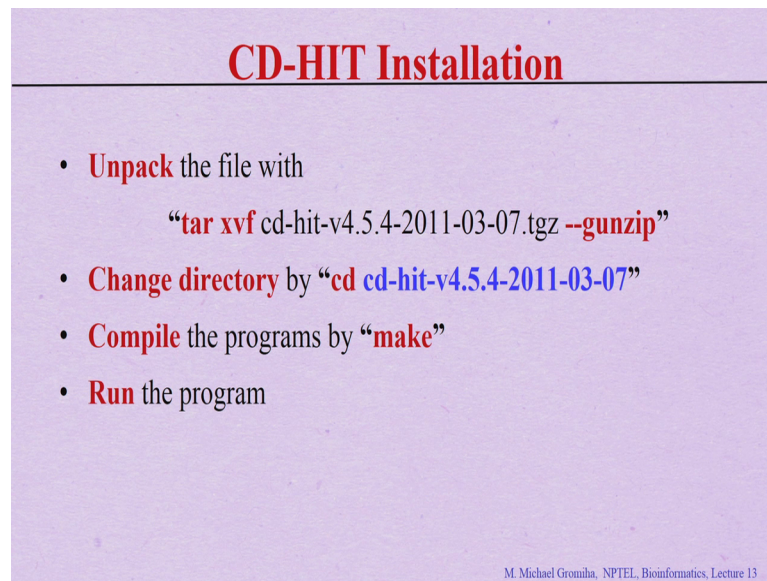
Student: 2

2, 3, 4, 5, 6, 7 times; 0.01 into 7; so, 0.11; so 0.11 is equal to 0.3 something. So, then this is the closest one? Here also you can find distance you can see this is the closest among these three sequences. So, you can make any pairs and you can make any different clusters, you can do any cutoff; I can do that. This here we used the amino acid

composition as a feature; likewise we can use any feature to assign the clusters and you can obtain the non redundant sequences.

So, how to install these CD HIT? This is freely available and this program is very easy to install and will work with this Linux operating systems.

(Refer Slide Time: 08:37)

A presentation slide with a light purple background. At the top, the title "CD-HIT Installation" is written in a bold, red, serif font. Below the title, there is a bulleted list of four steps for installation. Each step starts with a red bullet point and a red keyword: "Unpack", "Change directory", "Compile", and "Run". The instructions are as follows: 1. "Unpack the file with" followed by the command "tar xvf cd-hit-v4.5.4-2011-03-07.tgz --gunzip" in a red monospace font. 2. "Change directory by" followed by the command "cd cd-hit-v4.5.4-2011-03-07" in a red monospace font. 3. "Compile the programs by" followed by the command "make" in a red monospace font. 4. "Run the program". In the bottom right corner, there is a small, faint text credit: "M. Michael Gromiha, NPTEL, Bioinformatics, Lecture 13".

CD-HIT Installation

- **Unpack** the file with
`"tar xvf cd-hit-v4.5.4-2011-03-07.tgz --gunzip"`
- **Change directory** by `"cd cd-hit-v4.5.4-2011-03-07"`
- **Compile** the programs by `"make"`
- **Run** the program

M. Michael Gromiha, NPTEL, Bioinformatics, Lecture 13

So, first you download the files and then unpack using the tar and then you have to go to your directory. And you can install the program by using make, right and you can run the program. So, first it requires to download the program, you need to unpack and then to compile and you got to run the program. This is the steps we use when compelling the program right; in the UNIX systems.

(Refer Slide Time: 09:02)

Run CD-HIT

```
./cd-hit -i db -o db90 -c 0.9 -n 5
```

db: input file name 0.4 (40%)

db90: output file name

0.9, means 90% identity (clustering threshold)

5 is the size of word

M. Michael Gromiha, NPTEL, Bioinformatics, Lecture 13

So, how to run CD HIT? Because this is the program that you put dot slash CD HIT and this is the command we have to give. So, what is for i? That is input; we give this input file name; the db gives the input file name and what is o?

Student: Output.

This output; this is the output file name; db 90 is the output file name and here we give the c; 0.9 and n minus 5; what is 0.9?

Student: (Refer Time: 09:31)

This identity; so, which identity; you require whereas, 90 percent or 80 percent or 40 percent; if you need 40 percent; what you have to give?

Student: 0.4

0.4; for example this is 40 percent or n equal to 5; what is n equal to 5?

Student: Word size.

Word size; so in this case we look for the pentapeptides; how many times each pentapeptides occurs within those sequences? So, depending upon your sequence identity and the word size, so the CD HIT will run and then give you the sequences.

(Refer Slide Time: 10:04)

Run CD-HIT

Choice of word size:

- n 5 for thresholds 0.7 ~ 1.0
- n 4 for thresholds 0.6 ~ 0.7
- n 3 for thresholds 0.5 ~ 0.6
- n 2 for thresholds 0.4 ~ 0.5

M. Michael Gromiha, NPTEL, Bioinformatics, Lecture 13

For example n; this is a choice, they use the default say n equal to 5 for the thresholds 0.7 to 1; n equal to 4; they use 6 to 7 and 3 for 5 to 6; this is normal threshold we use for running the CD HIT.

(Refer Slide Time: 10:18)

Example

input *output*

`./cd-hit -i hemoglobin.fasta -o db85 -c 0.85 -n 5`

```
>sp|P69905|HBA_HUMAN Hemoglobin subunit alpha OS=Homo sapiens GN=HBA1 PE=1 SV=2
MVLSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSPPTTKTYFPHFDLSHSAQVKGHG
KVVADALTNAAVHVDMPNHALSALSDLAHMLRVDPVNFKLLSHCLLVTLAAHLPAETP
AVHASLQKFLASVSTLTSKYR
>sp|P01945|HBA_JAT Hemoglobin subunit alpha-1/2 OS=Battus norvegicus GN=HBA1 PE=1 SV=3
MVLSDQKTNKIKWQKIGGHGGEYGEALQRMFAAPPTTKTYFPHFDLSHSAQVKGHG
KVVADALNAADVELDQALSTLSDLAHMLRVDPVNFKLLSHCLLVTLAAHHPDFT
AVHASLQKFLASVSTLTSKYR
>sp|P01941|HBA_MOUSE Hemoglobin subunit alpha OS=Mus musculus GN=HBA1 PE=1 SV=2
MVLSDQKTNKIKWQKIGGHGGEYGEALQRMFAAPPTTKTYFPHFDLSHSAQVKGHG
KVVADALASAQHLSDLPALSDLAHMLRVDPVNFKLLSHCLLVTLAAHHPDFT
AVHASLQKFLASVSTLTSKYR
>sp|P01966|HBA_BOVIN Hemoglobin subunit alpha OS=Bos taurus GN=HBA1 PE=1 SV=2
MVLSDQKTNKIKWQKIGGHGGEYGEALQRMFAAPPTTKTYFPHFDLSHSAQVKGHG
KVVADALNAADVELDQALSTLSDLAHMLRVDPVNFKLLSHCLLVTLAAHHPDFT
AVHASLQKFLASVSTLTSKYR
>sp|P01958|HBA_HORSE Hemoglobin subunit alpha OS=Equus caballus GN=HBA1 PE=1 SV=2
MVLSDQKTNKIKWQKIGGHGGEYGEALQRMFAAPPTTKTYFPHFDLSHSAQVKGHG
KVVADALNAADVELDQALSTLSDLAHMLRVDPVNFKLLSHCLLVTLAAHHPDFT
AVHASLQKFLASVSTLTSKYR
>sp|P01997|HBA_ZANTR Hemoglobin subunit alpha OS=Pan troglodytes GN=HBA1 PE=1 SV=2
MVLSDQKTNKIKWQKIGGHGGEYGEALQRMFAAPPTTKTYFPHFDLSHSAQVKGHG
KVVADALNAADVELDQALSTLSDLAHMLRVDPVNFKLLSHCLLVTLAAHHPDFT
AVHASLQKFLASVSTLTSKYR
>sp|P01959|HBA_EQUUS Hemoglobin subunit alpha OS=Equus asinus GN=HBA1 PE=1 SV=3
MVLSDQKTNKIKWQKIGGHGGEYGEALQRMFAAPPTTKTYFPHFDLSHSAQVKGHG
KVVADALNAADVELDQALSTLSDLAHMLRVDPVNFKLLSHCLLVTLAAHHPDFT
AVHASLQKFLASVSTLTSKYR
>sp|P01963|HBA_PIG Hemoglobin subunit alpha OS=Sus scrofa GN=HBA1 PE=1 SV=1
MVLSDQKTNKIKWQKIGGHGGEYGEALQRMFAAPPTTKTYFPHFDLSHSAQVKGHG
KVVADALNAADVELDQALSTLSDLAHMLRVDPVNFKLLSHCLLVTLAAHHPDFT
AVHASLQKFLASVSTLTSKYR
>sp|P06633|HBA_PONNY Hemoglobin subunit alpha OS=Pongo pygmaeus GN=HBA1 PE=2 SV=2
MVLSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSPPTTKTYFPHFDLSHSAQVKGHG
KVVADALTNAAVHVDMPNHALSALSDLAHMLRVDPVNFKLLSHCLLVTLAAHLPAETP
AVHASLQKFLASVSTLTSKYR
>sp|P06739|HBA_CANIS Hemoglobin subunit alpha OS=Canis familiaris GN=HBA1 PE=1 SV=1
MVLSPADKTNKIKWQKIGGHGGEYGEALQRMFAAPPTTKTYFPHFDLSHSAQVKGHG
KVVADALNAADVELDQALSTLSDLAHMLRVDPVNFKLLSHCLLVTLAAHHPDFT
AVHASLQKFLASVSTLTSKYR
```

M. Michael Gromiha, NPTEL, Bioinformatics, Lecture 13

So, now we give the sequences; so, we have all the sequence in fasta format and we run CD HIT; this is the input file, and here this is the output. So, with 85 percent sequence identity and using the word size n; so, you can use any word size right here n equal to 5. So, then you see you will get the final data.

(Refer Slide Time: 10:40)

Example

```
[gromiha@INSIGHT1 cd-hit-v4.5.4-2011-03-07]$ more db85
>sp|P69905|HBA_HUMAN Hemoglobin subunit alpha OS=Homo sapiens GN=HBA1 PE=1 SV=2
MYLSPADKTNVKAANGKVGAGAGETGAEALQRMFAAFFTTKTYFPHPDLSHGSAQVKHG
KKVADALTNVAHVDDMPNALSALSDLHAHKLKRVDPVNFKLLSHCLLVTLAAHLPAEFTF
AVHASLQKFLASVSTVLTSKYR
>sp|P01946|HBA_RAT Hemoglobin subunit alpha-1/2 OS=Rattus norvegicus GN=Hba1 PE=1 SV=3
MYLSADDKTNIKNCWKGKGGHGETGEEALQRMFAAFFTTKTYFPHPDLSHGSAQVKAHG
KKVADALAKAAADHVDELPGALSTLSDLHAHKLKRVDPVNFKLLSHCLLVTLACHHPDFTF
AMHASLQKFLASVSTVLTSKYR
>sp|P01965|HBA_PIG Hemoglobin subunit alpha OS=Sus scrofa GN=HBA PE=1 SV=1
VLSAADKAMVKAANGKVGQAGAGAEALRMFLGFTTKTYFPHPDLSHGSDQVKAHQ
KVADALTKAVGHLDLPGALSALSDLHAHKLKRVDPVNFKLLSHCLLVTLAAHHPDDPNP
VHASLQKFLANVSTVLTSKYR
>sp|P60529|HBA_CANFA Hemoglobin subunit alpha OS=Canis familiaris GN=HBA PE=1 SV=1
VLSPADKTNIKSTWOKGGHAGDYGGEALDRFTQSFPTTKTYFPHPDLSHGSAQVKAHGK
KVADALTTAVAHLDLPGALSALSDLHAYKLKRVDPVNFKLLSHCLLVTLACHHPTEPTPA
VHASLQKFFAAVSTVLTSKYR
[gromiha@INSIGHT1 cd-hit-v4.5.4-2011-03-07]$
```

Connected to 10.93.219.140 SSH2 - aes12

M. Michael Gromiha, NPTEL, Bioinformatics, Lecture 13

Now, you will get the data then you can use this data for any analysis. And now you have two sets of data; one is the redundant sequence, one is non redundant sequence. If you compare your composition or your any properties; you can see the difference between the non redundant ones as well as redundant ones.

Also it depends upon; how much your data is reduced. If you can see much reduction in your data sets, then you can see much deviation from the properties; what you calculate from the whole data set or from the reduced data set. There is not much redundancy, that already your data set is non redundant; then you cannot see much difference. Then we will get similar data; when you analyze with a whole set of data or with the reduced set of data.

(Refer Slide Time: 11:27)

Blastclust

- **Blastclust** is a program within the standalone BLAST package used to cluster either protein or nucleotide sequences.
- The program begins with pairwise matches and places a sequence in a cluster if the sequence matches at least one sequence already in the cluster.
- In the case of proteins, the blastp algorithm is used to compute the pairwise matches.

M. Michael Gromiha, NPTEL, Bioinformatics, Lecture 13

So, there is another program; this is also widely used program one of the applications of Blastclust is you can get the data with any redundancy. Even if you want to have 20 percent identity or 25 percent identity; you can use with this Blastclust. This is the standalone program which is available within BLAST and you can use this to cluster either protein sequence or nucleotide sequence. You can use it for both protein as well as the DNA sequences to get the clusters.

So, how it works? It must begin with the pairwise match and then places the sequence in a cluster. If the sequence matches at least one with the other, if two sequences have high identity; then if you are interested to threshold at 25 percent; which is more than 25 percent, you put one in the clusters. In the case of the proteins, they used a blastp for the alignment; for the case of nucleic acids, which program they use?

Student: blastn

blastn. So, you take the blastn to align the sequences.

(Refer Slide Time: 12:31)

Blastclust

The general command to create a set of non-redundant set of protein sequences is

```
blastclust -i infile -o outfile -p T -L .9 -b T -S 95, 7A
```

where infile and outfile are input and output files, respectively.

T stands for protein;

The coverage of the length and sequence identity cutoff are 90% (-L .9) and 95% (-S 95), respectively.

AITLVIKS
ITVAIKS
Coverage: 70%

M. Michael Gromiha, NPTEL, Bioinformatics, Lecture 13

Then you get, this is a command; which you can use to extract the data. You can first get the program using the standalone program, you can get non redundant sequences using this sequence. So here it's the same; i is the input file, so infile is the input file and outfile is the output file; then we get some numbers -p T; -L, .9; -b T, -S 95.

So, what this T stands for? Protein, for protein say put T. For example, they have the L; L equal to 0.9 and what is a meaning of L equal to 0.9?

Student: coverage

It is a coverage; so, its coverage is about, that is 90 percent; what is the meaning of coverage?

Student: Sequence coverage.

Sequence coverage you know for example, if we have one sequence; this is AITLVIKS; now you can see aligned only some part of the sequences; for example, aligned only here, in this case what is the query coverage?

Student: 80? 70 percent.

Coverage equal to?

Student: 70 percent.

70 percent; so, this is how much coverage you can accept. And then you can see this 95; what is S 95?

Student: Identity.

Identity; you can say this about 95 percent sequence identity. So, here you can restrict with a sequence identity plus the coverage; what is the advantages of using coverage?

Student: how much you have.

For example, if you have 1 proteins; 100 residues, another proteins 200 residues. Only some part; only 10 residues are aligned and 10 residues are same. If you do not use your query coverage, the aligned 10 sequences and 10 sequences are same. So, what will be the identity?

Student: 100.

100 percentage, but you see the query coverage out of 100; with aligned only 10. So, coverage only 10 percent and that means; all the 90 are totally different. In this case, they are different from each other; so to include this information; they use the query coverage in this program. Right now, we discussed about two different algorithms or the two different programs.

Student: CD HIT.

CD HIT.

Student: Blastclust.

(Refer Slide Time: 14:49)

PISCES

- PISCES is a protein sequence culling server to produce subsets of non-redundant sequences using Protein Data Bank entries or Uniprot sequences in FASTA format.
- Sequence identities for PDB sequences are determined by the combination of Combinatorial Extension **structural alignment** and **PSI-BLAST alignment**.

M. Michael Gromiha, NPTEL, Bioinformatics, Lecture 13

And the blastclust; so, there is another program that's called PISCES. This is a protein sequence culling server; to get the subsets of non redundant sequences. Here we have different options; either you can use the protein sequences already available in this server; mainly from the PDB. Or you can give your own sequence; for example, if you only give your own sequence, you can give the UniProt sequences in fasta format. If you want to compare with the existing ones; then you can use the existing PDB data and you can compare and you can get any non redundant sequences.

So, how they get this non redundancy? If you have the PDB sequence; they use a structural alignment and the PSI-BLAST alignment. So what is PSI-BLAST alignment? It is a kind of multiple sequence alignment, there is a position; specific iterated BLAST alignment. So, they have the multiple sequence alignment by putting same sequences or similar sequences at the first. So, they have this sequence alignment as well as a structured based sequence alignment; they are combined together, if you have the PDB entries.

(Refer Slide Time: 15:54)

PISCES

- Non-PDB sequences are culled with **sequence identities from PSI-BLAST**. PISCES does not search the non-redundant sequence database, but rather use the user's input sequences as the database.
- This server will usually be used to cull a related set of sequences, for instance those from a PSI-BLAST search.
- It takes the amino acid sequence in FASTA format and sends the list of non-redundant protein sequences by e-mail.

→ <http://dunbrack.fccc.edu/pisces/>

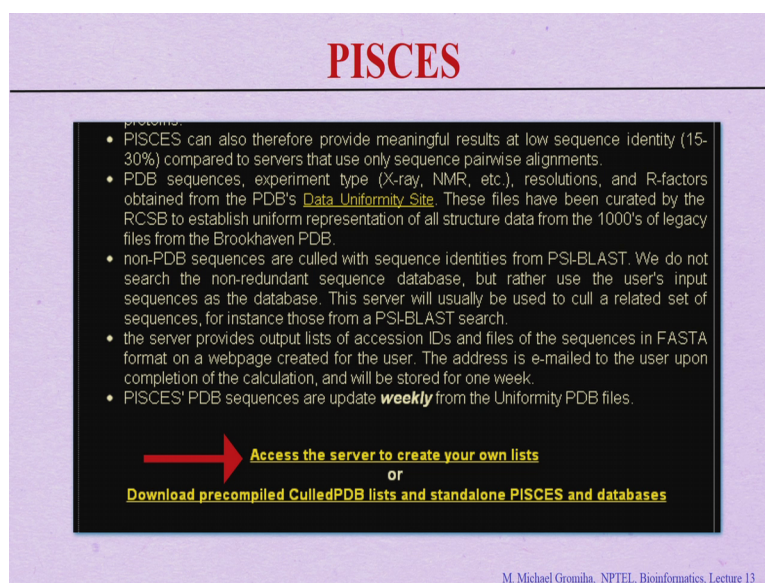
M. Michael Gromiha, NPTEL, Bioinformatics, Lecture 13

If they do not have this PDB entries; for example, non PDB; then they get this sequence identity only from PSI-BLAST; it is because structures are not known. So, they have these PSI-BLAST sequences, they do the homologues sequences and this will give you the values and the non redundant sequences using only the sequence information.

So, this you can use this to cull any protein sequences; I did various ranges of sequence identities. So, what is the input sequence for this PISCES server? It takes amino acid sequence in fasta format as the input and then give the non redundant sequences. So, this is the web server; so where you can submit your query sequences and you will get the output either by email or you can get download.

If there is less number of data; you can get it immediately, but the problem with this file server is it can handle only less number of sequences. There is some limitations about the query sequence, when in the case of CD HIT or theBlastclust; you can get any number of sequences; it can handle huge amount of sequences. So, that is a difference between the CD HIT as well as with the PISCES servers.

(Refer Slide Time: 17:03)



PISCES

proteins.

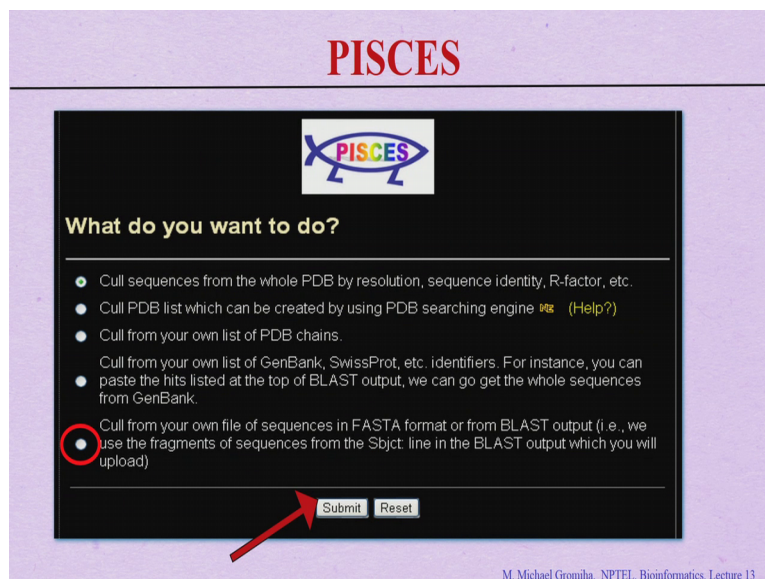
- PISCES can also therefore provide meaningful results at low sequence identity (15-30%) compared to servers that use only sequence pairwise alignments.
- PDB sequences, experiment type (X-ray, NMR, etc.), resolutions, and R-factors obtained from the PDB's [Data Uniformity Site](#). These files have been curated by the RCSB to establish uniform representation of all structure data from the 1000's of legacy files from the Brookhaven PDB.
- non-PDB sequences are culled with sequence identities from PSI-BLAST. We do not search the non-redundant sequence database, but rather use the user's input sequences as the database. This server will usually be used to cull a related set of sequences, for instance those from a PSI-BLAST search.
- the server provides output lists of accession IDs and files of the sequences in FASTA format on a webpage created for the user. The address is e-mailed to the user upon completion of the calculation, and will be stored for one week.
- PISCES' PDB sequences are update **weekly** from the Uniformity PDB files.

→ [Access the server to create your own lists](#)
or
[Download precompiled CulledPDB lists and standalone PISCES and databases](#)


M. Michael Gromiha, NPTEL, Bioinformatics, Lecture 13

So, this is the online tool; so if you want to get the sequences which are already available in the server; you can get it. Or if you want to create your own data set; so, for example if you have 1000 sequences or you want to get the non redundant sequences; in this case you can get this, access the server; to create our own list.

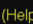
(Refer Slide Time: 17:26)



PISCES



What do you want to do?

- Cull sequences from the whole PDB by resolution, sequence identity, R-factor, etc.
- Cull PDB list which can be created by using PDB searching engine  ([Help?](#))
- Cull from your own list of PDB chains.
- Cull from your own list of GenBank, SwissProt, etc. identifiers. For instance, you can paste the hits listed at the top of BLAST output, we can go get the whole sequences from GenBank.
- ☒ Cull from your own file of sequences in FASTA format or from BLAST output (i.e., we use the fragments of sequences from the Sbjct: line in the BLAST output which you will upload)

Submit Reset

M. Michael Gromiha, NPTEL, Bioinformatics, Lecture 13

If you click on that; this will ask for the different options, which one you want. So, you can if you have your own sequence; you can keep the last one; here you can see your own file of sequences then if you submit.

(Refer Slide Time: 17:39)

The screenshot shows the PISCES web interface. At the top, the title 'PISCES' is displayed in red. Below it, a subtitle reads '>>PISCES --server: Taking input parameters for culling protein sequences'. The main content area has a black background with white text. It prompts the user to 'Please browse or paste FASTA format sequences or BLAST/PSI-BLAST output file'. There is a text box for pasting sequences, containing a sample FASTA entry for a protein. Below the text box is an 'Upload file:' label and a 'Browse...' button. Further down, the 'Set sequence identity threshold:' section contains three input fields: 'Maximum percentage identity' (set to 25), 'Minimum chain length' (set to 40), and 'Maximum chain length' (set to 10000). At the bottom of this section are 'Submit' and 'Reset' buttons. A red arrow points to the 'Submit' button. The footer of the slide includes the text 'M. Michael Gromiha, NPTEL, Bioinformatics, Lecture 13'.

Then you can get this next page; here you have to paste or type your fasta sequence in this text box. You have the sequences; you can give the text box here or you can upload file, they have the upload file option; you can browse from your computer and you can upload the file which contains amino acid sequences in which format?

Student: Fasta.

Fasta format because it accepts fasta format of any UniProt sequences. Then you can also give your options; what is the maximum percent identity you want? 25 percent or 30 percent or 40 percent or right; so, you can give the threshold as per your convenience; depending upon the problem, depending upon initial data sets and your requirement; you can set the threshold.

Normally we get 25 percent as a threshold; then they asking about the chain length, whether we need to consider all the sequences or discard some short sequence and large sequence. So, here we can give 40 or 50 because several proteins are more than 50 residues; only few exceptions. So, in this case you can give the minimum length; why do you needing the minimum length? This can avoid considering the peptides; there are many short peptides. In this case you can avoid these peptides from these not redundant data set and maximum chain length that also you can specify; for your calculation purpose; then it takes all these information as input.

(Refer Slide Time: 19:08)

Sequence ID		length
> 7467903 (Genbank) Outbrank membrane		
MKFTTITITITFTLSTYVIALDLQLALGTQYKNSKGLAKQKIFLALIE	Id#	456
QGFQKSGFNGVLIGQIMQKRTKCYVYKGLVSTFSTAGTGLITL	11559475 (Genbank) Outbrank	783
EGFLFQKSDALLAQGLQYKRSFSTAGTQYKGLVSTFSTAGTGLITL	14695241 (Genbank) Outbrank	785
STLTLSESPVTHLQVYVTEKSLGGLDILDAASAGLQKGLVSTFSTAGTGLITL	15569606 (Genbank) Outbrank	256
ATAGCGKANEIVTGLSDILTMQGLFRLIQLSDEGTAKSAFNPDI	P24061 (Genbank) Outbrank	452
NARHIVYITLALQGLVSTFSTAGTGLITL	15597487 (Genbank) Outbrank	452
YTTLITLITITITFEGKAGLQYKRSFSTAGTQYKGLVSTFSTAGTGLITL	15597487 (Genbank) Outbrank	452
WQSFSTAFVIAVQGLVSTFSTAGTGLITL	15597487 (Genbank) Outbrank	452
EKMTITQYKNSVIAVQGLVSTFSTAGTGLITL	15597487 (Genbank) Outbrank	452
PEITG	> 11559475 (Genbank) Outbrank integral membrane	
MTNRPNRPITATLVLGIVISVQAQVYVYKGLVSTFSTAGTGLITL	11559475 (Genbank) Outbrank	783
ATITLITFVQVQGLVSTFSTAGTGLITL	14695241 (Genbank) Outbrank	785
VYFEGADLQVPTTIFQKQFQVQVYVYKGLVSTFSTAGTGLITL	P13949 (Genbank) Outbrank	201
SPNTTATVYPAINTVYVYKGLVSTFSTAGTGLITL	161945 (Genbank) Outbrank	292
NALDLAGLQVQGLVSTFSTAGTGLITL	7208425 (Genbank) Outbrank	560
ATGLVQGLVSTFSTAGTGLITL	15596061 (Genbank) Outbrank	260
SEINAFNSWQGLVSTFSTAGTGLITL	1347083 (Genbank) Outbrank	794
QGLVSTFSTAGTGLITL	P19196 (Genbank) Outbrank	935
YVYKGLVSTFSTAGTGLITL	12620510 (Genbank) Outbrank	201
QGLVSTFSTAGTGLITL	P16001 (Genbank) Outbrank	930
QGLVSTFSTAGTGLITL	15596468 (Genbank) Outbrank	616
YVYKGLVSTFSTAGTGLITL	P15727 (Genbank) Outbrank	482
YVYKGLVSTFSTAGTGLITL	P20765 (Genbank) Outbrank	778
YVYKGLVSTFSTAGTGLITL	P14666 (Genbank) Outbrank	1577
YVYKGLVSTFSTAGTGLITL	3228547 (Genbank) Outbrank	700
YVYKGLVSTFSTAGTGLITL	P06970 (Genbank) Outbrank	812
YVYKGLVSTFSTAGTGLITL	P22340 (Genbank) Outbrank	505
YVYKGLVSTFSTAGTGLITL	P44601 (Genbank) Outbrank	565
YVYKGLVSTFSTAGTGLITL	P35077 (Genbank) Outbrank	584
YVYKGLVSTFSTAGTGLITL	17212580 (Genbank) Outbrank	444
YVYKGLVSTFSTAGTGLITL	7474849 (Genbank) Outbrank	654
YVYKGLVSTFSTAGTGLITL	12913794 (Genbank) Outbrank	590
YVYKGLVSTFSTAGTGLITL	P06111 (Genbank) Outbrank	257
YVYKGLVSTFSTAGTGLITL	P41265 (Genbank) Outbrank	530
YVYKGLVSTFSTAGTGLITL	P39041 (Genbank) Outbrank	719
YVYKGLVSTFSTAGTGLITL	1595231 (Genbank) Outbrank	462

M. Michael Gromiha, NPTEL, Bioinformatics, Lecture 13

And if you submit then this will give you the sequence, here this is the complete sequence and this is the list for the ID's. You can see the id here and here we give the complete sequences and here we can give sequence ids.

So, you can use a PISCES to get the non redundant sequences, but the only disadvantage of PISCES is; it can handle only less number of sequences then it uses a clustering technique, then they will finally, they give you the final set of non redundant sequences. So, till now what did we discuss in the class? today's lecture?

Student: Redundancy.

First we discussed about the redundancy, so what is redundancy?

Student: Over representation of sequences

Right we can see if there are same sequences; if one or more same sequences which are present in any data set. If we construct in a specific data set which consists of 100 sequences; are there any sequences or any 2 or 3 sequences which share the common information with the high sequence identity. So, if we share high sequence identity then these sequences are called redundant sequences; then how to get the non redundant sequences.

Student: (Refer Time: 20:29)

For example, if you have two sequences; different sequence identity and a third one this has the different sequence identity. If you want to have the sequences with more than 50 percent then what we have to do?

Student: (Refer Time: 20:43)

When the first two it is very high; so you can keep one and discard another one. Because if you have the same sequences that will reduce a bias in the final calculations as well as for any prediction algorithms. We checked two examples like what will happen if the same sequences are present two times or three times, depending upon the features; some of them are over represented and some of them are underrepresented, should not avoid the bias; we need to construct non redundant sequences.

So, what are the various program we discussed to construct the non redundant sequences?

Student: CD HIT

CD HIT.

Student: Blastclust

Blastclust.

Student: PISCES

And the PISCES; so, what is the advantage of using CD HIT?

Student: It can handle large dataset.

It can have the large data sets.

Student: very fast

It is very fast; so, the disadvantages of this CD HIT.

Student: upto 30 percent

It can handle up to 30 to 40 percent sequence identity because it has two versions; one is online and one is the standalone version. So, you can handle only up to 30 percent or 40

percent sequence identity based on the program you use. So, what principle it used for clustering the sequences?

Student: Hamming distance

K means clustering; to make these clusters, how this K means clustering works? Try to consider all the sequences and we can do some number of clusters. So, this K clusters and then try to put in the sequences, similar sequences within that particular cluster and they rearrange all the values again and again. So, that the minimum redundancy is maintained within that only particular cluster.

This would take any cluster, you can see there is completely close to each other then they take the representatives from different clusters; in this case you can get the minimum redundancy among your sequences. So, what are various features we use for clustering techniques? We discussed two aspects; both are based on?

Student: distance.

This amino acid composition; so, one is the Hamming distance and the Euclidean distance. So, if you have any pair of sequences; you can calculate the values and make any threshold, if the pairs are more than that threshold then you can discard one. So, in getting the non redundant sequences; some advantages, some disadvantages. For example, if we have 100 sequences; made into 10 clusters, first cluster has 20, second cluster has 5 and another cluster has 1. Then how to get the representative data; if it is 1; it is easy, you can then pick it up; if that is 10 how to you get the representative one?

Student: the common one

Generally, we will take randomly pick up one; you can get this one. This is the reason; you have to cluster again and again and you can see what will happen if you pick up these different sequences from each cluster. You can see almost they are in similar features; so you can get similar data, even then you can do this test to confirm that your data are significant or not. So, then the second one is a Blastclust; what is the advantages of Blastclust?

Student: It can also handle large data

It can also handle large data sets and you can go with the less sequence redundancy; even 20 percent, 25 percent you can do. And it also considers the coverage; query coverage, so with that you can use the Blastclust for clustering the sequences. The third one we discussed about PISCES; here it is easy to use, because if you give all the sequences on length; you will get the non redundant sequences.

But only disadvantage is; it can handle limited number of sequences. In this case, you have to do this again and again; take your sequence and get non redundancy and then make another set and get the non redundant sequences and cluster together and put again, you have to run again to get the non redundant sequences.

So, we discussed about the construction of non redundant sets and various features of the primary sequence. There are various parameters we calculated and we used for different applications. Now in the next class, we will discuss about the next level; what is next level of protein structure? Secondary structures; what are different secondary structures?

Student: alpha helices

Alpha helices beta strands coil and all; so, we will discuss about what are the different secondary structures? What are the probable allowed positions of these secondary structures in the construction of Ramachandran plot? How to predict these secondary structures from any given amino acid sequence? What is the reliability of these prediction techniques and so on; that we will discuss in the next following couple of classes.

Thank you for your kind attention.