

Assignment 2 / Biometry and bioinformatics I / 2014

Monitoring/tracing of viruses and bacteria on the basis of their informative sequences is done very actively, providing understanding of their behaviour, for example spread of an epidemic. In assignment set 4 you will get familiar with various kind of examples by reading and writing short essay(s) or comments about scientific publications.

In this assignment 2 you get familiar with one database (out of many different kind of virus- and bacteria sequence databases) and use it's information (seqs) for studying the behavior of the influenza-virus H1N1 which caused a pandemia in 2009.

- The database includes extensive sequence information. What does this mean? If a database has, say, 100 seqs from the virus from a given country from a given time point (year, for example), this does not mean that there were 100 infections. It means that 100 sequences have determined, according to some criteria, according to some interests within the scope of research financing, for example. One seq (with an accession number, is a virus from one infected human (or animal) individual.

Database: <http://www.ncbi.nlm.nih.gov/genomes/FLU/FLU.html>

The study question: on the basis of existing sequence information, what can be inferred about the behaviour of H1N1 by sequence clustering

- before the pandemia -> 2008
- during the pandemia 2009 and after that (2010 ->)
- The most active period of the pandemia was Autumn 2009. Your task is to first pick up countries which you would like to study, and from which information exist, perform an initial study (part (a)) and then a more comprehensive study (part (b)).
- The pandemia started from Mexico in Spring 2009. There is lots of data from various countries from 2009 and also 2010 ->. But very little data from the period before 2009. You don't have to restrict your study so that all your countries should have pre-pandemic data! Take at least one country from which there is pre-pandemic data, and data from 2009 and 2010 ->. And then additional countries with data from 2009 and 2010 ->.
- Demo of the flu-database during the first course session (9.9) => starting data collection.
- At the beginning of the second course session a short demo about MEGA5, the software for sequence clustering => clustering by neighbor joining algorithm can be done after this demo; a more detailed lecture about this kind of sequence processing methods will be in session 23.9.

Part (a)

First you perform a small experiment by collecting a small set of seqs from 2-3 countries, from different time points. Small set = 1-2 seqs as representatives of one country and one time point. Align the seqs and cluster them by the neighbor joining method.

- You get some kind of an idea to perform a more comprehensive study.
- In your final report give also this initial neighbor joining tree and explain your initial interpretation.
- For example "it seems to be that clustering (sequence similarities – differences) is based on (x) geographical location, (y) time; you very probably notice that either (x) or (y) is the predominant "apparent explanation". Describe your explanation (i.e. "the seems to be" – situation)

Part (b)

On the basis of your “seems to be” –experiment (part (a) you make a plan for a “real study” (part (b)) by collecting more data.

The amount of data is your own choice. Maybe more seqs (now some tens; in part (a) you had one or two) from those countries which you included in your experiment, maybe other countries (nearby vs. remote), maybe more time points.

Then you again align the seqs, perform neighbor joining clustering and write your interpretation about the epidemic behavior of H1N1.

Your report should thus include the neighbor joining -tree with explanations from part (a) and neighbor joining -tree with explanations from part (b). And the study plan which you made, on the basis of part (a) to do part (b). A study plan very probably includes modifications to your original plan because of lack of data.

Practical advise

H1N1 is one of many different influenza A-virus subtypes.

Influenza A-virus is composed of 11 genes, H (HA) and N (NA) being those genes which serve as defining a given subtype.

H = hemagglutinin gene (HA)

N = neuraminidase gene (NA)

When a virus is typed by sequencing, usually at least the H-gene is sequenced => most information is from the H-gene.

We restrict this course assignment to H-gene sequences from H1N1 and we work with DNA-information (nucleotide information), not protein (amino acid information).

Inspect carefully the contents of next page: how to find H-gene nucleotide seqs, i.e. what boxes you should click in the database in order to define that you want to get H-seqs (you must click HA-protein) from the H1N1 type of A-influenza virus (you must click 1 from the H-subtype box and 1 from the N-subtype box). And, when collecting data, you should download “nucleotide”, not “protein”.

Practical advice for collecting data from Flu-database

<http://www.ncbi.nlm.nih.gov/genomes/FLU/FLU.html>

The main page has links to all kind of updated information about flu-viruses

Go here



The screenshot shows the top of the InFluenza Virus Resource website. It features the NCBI logo on the left, the title "InFluenza Virus Resource" with the subtitle "Information, Search and Analysis" in the center, and a colorful virus icon on the right. Below the title is a navigation menu with links: Home, Search, Reference, Flu home, Database, Genome Set, Alignment, Tree, BLAST, Annotation, FTP, Help, and Contact us. A descriptive paragraph follows, stating that the resource presents data from the NIAID Influenza Genome Sequencing Project and GenBank, along with tools for sequence analysis and submission. At the bottom, there are links to "Read more about" various aspects of the resource.

(The database has many facilities. Such as aligning data, drawing phylogenetic trees. You can, of course, try them. However, they do not operate very well.....)

The screenshot shows the "Define search set" interface of the InFluenza Virus Resource. It includes several dropdown menus for filtering search results: Type (A, B, C), Host (Feline, Ferret, Giant anteater, Human), Country/Region (any, regions, Northern temperate, Southern temperate), Protein (PA, PA-X, P3, HA), Subtype (H1, H2, H3, N1, N2, N3), and Sequence length (Full-length only, Full-length plus). There are also fields for Collection date and Release date, and a checkbox for "collapse identical sequences".

Influenza A-virus (with an arrow pointing to the Type dropdown)

Select from here countries/regions you want to study. (with an arrow pointing to the Country/Region dropdown)

Read this . Not a good idea to collect many identical copies. On the other hand, which seqs are present as many copies, is valuable information. (with an arrow pointing to the "collapse identical sequences" checkbox)

These three clicked on specify that you will get hemagglutinin seqs from the type H1N1-virus (with arrows pointing to the Protein dropdown set to "HA", the Subtype dropdown set to "H1", and the N dropdown set to "N1")

2009-2013. Lots of data from 2009, some data may exist also 2010 -> (with an arrow pointing to the Collection date fields)

An example, results from Finland, you need nucleotide data, click from here
 From here you can define the information appearing in the title of each seq

It is not a good idea to keep the title (= text appearing after >... too long). Include country, year, month, accession number.

The screenshot shows the NCBI Influenza Virus Resource interface. At the top, there are navigation tabs: 'File home', 'Database', 'Genome Set', 'Alignment', 'Tree', 'BLAST', 'Annotation', 'Submission', 'FTP', and 'Virus resources'. Below these are buttons for 'Add your own sequences', 'Do multiple alignment', 'Build a tree', 'Download', and a dropdown menu currently set to 'Nucleotide (FASTA)'. A 'Define' field is visible with a dropdown arrow, and a 'Customize FASTA define' link. Below the search options, there is a 'Show query' link and a note: 'Note: All groups of identical sequences in the dataset will be represented by the oldest sequence in the group. Metadata of the collapsed sequences are not preserved.' The main content is a table titled '74 protein sequences after collapsing (168 total)'. The table has columns: Accession, Length, Host, Protein, Subtype, Country, Region, Date, Virus name, Mutations, Age, Gender, Lineage, VacStr, Complete, and #. The first few rows are highlighted in blue.

| Accession | Length | Host | Protein | Subtype | Country | Region | Date | Virus name | Mutations | Age | Gender | Lineage | VacStr | Complete | # |
|--|--------|-------|---------|---------|---------|--------|------------|---|-----------|-----|--------|---------|--------|----------|---|
| <input checked="" type="checkbox"/> AEN88650 | 566 | Human | HA | H1N1 | Finland | N | 2010/12/07 | Influenza A virus (A/Finland/15/2010(H1N1)) | | | | | | c | |
| <input checked="" type="checkbox"/> AEN88651 | 566 | Human | HA | H1N1 | Finland | N | 2010/12/09 | Influenza A virus (A/Finland/17/2010(H1N1)) | | | | | | c | 3 |
| <input checked="" type="checkbox"/> AEN88652 | 566 | Human | HA | H1N1 | Finland | N | 2010/12/09 | Influenza A virus (A/Finland/19/2010(H1N1)) | | | | | | c | 3 |
| <input checked="" type="checkbox"/> ADM95864 | 566 | Human | HA | H1N1 | Finland | N | 2010/01/12 | Influenza A virus (A/Finland/2/2010(H1N1)) | | | | | | c | |
| <input checked="" type="checkbox"/> AEN88656 | 566 | Human | HA | H1N1 | Finland | N | 2010/12/09 | Influenza A virus (A/Finland/22/2010(H1N1)) | | | | | | c | 3 |
| <input checked="" type="checkbox"/> AEN88651 | 566 | Human | HA | H1N1 | Finland | N | 2010/12/09 | Influenza A virus (A/Finland/27/2010(H1N1)) | | | | | | c | |
| <input checked="" type="checkbox"/> ADM95865 | 566 | Human | HA | H1N1 | Finland | N | 2010/01/13 | Influenza A virus (A/Finland/3/2010(H1N1)) | | | | | | c | |

Aligning seqs:

- You can use either Clustal (installed in classroom C128 computers; freely downloadable to your own computer, too) or the MAFFT-server.
- In both cases you take the FASTA-file, now aligned, from the result-link and start operating with it.
- First you must delete the extra / unnecessary and confusing part from the beginning (if such exists): You can easily see that while most seqs start from the *real start position of the gene*, ATG..., some seqs have something extra before ATG. Check how many nucleotides (in those cases where some extra nucleotides exist, there are usually up to 46 bp extra, check !). Remove the extra block and save the edited (= extra nucleotides removed) file as a new FASTA-file and work further with that file.

Report from Assignment 2 to be submitted to course Moodle during the first week of October. Your report should also include your aligned FASTA-files as separate text documents. In your report document: explain what you done and your interpretations about the behaviour of H1H1 on the basis of your study scheme. Copy-paste your neighbor joining clustering tree(s) into your report (DO NOT include them as separate MEGA5-files.)

During the last course session we shall see look together your work and also in the light of scientific literature about H1H1.

One of the extra assignments, for extra credits (i.e. extending the 5 cr -> more credits to added later without a deadline) will a more comprehensive study on H1N1 also including scientific literature as a baseline.