# BIOINFORMATICS SESSION 11. PRACTICE

2023-12-4
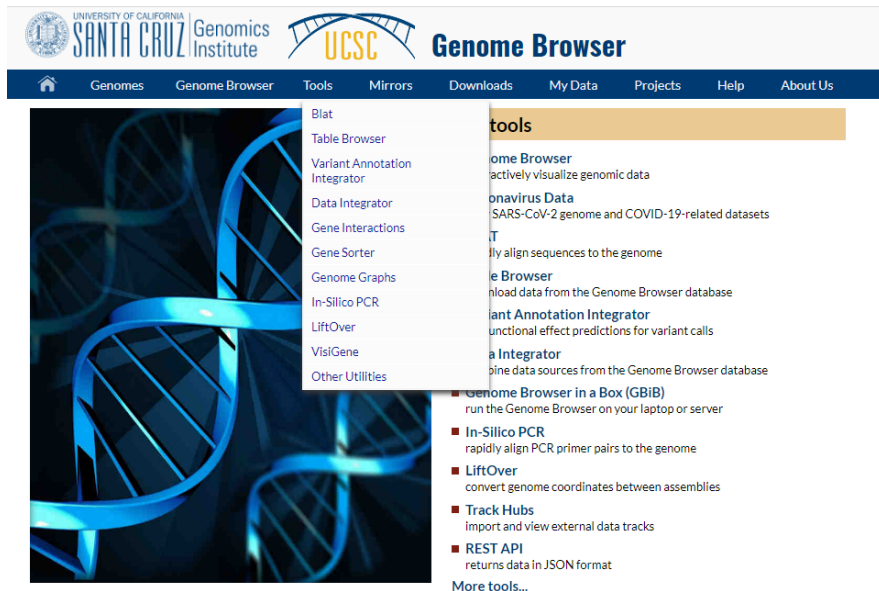
Personal genomes:
The differences between you and me

# Counting SNPs

- Using the table browser at the UCSC genome database
  (http://genome.ucsc.edu/cgi-bin/hgTables?org=human)

- Comparing chr4 of eight different human individuals
  (1) YanHuang (Han Chinese individual, anonymous)
  (2) Seong-Jin Kim (Korean)
  (3) James Watson
  (4) Craig Venter
  (5) YRI NA18507 (Yoruba, anonymous of the 1000 Genomes Project)
  (6) NA12891 (Central European origin, anonymous of the 1000 Genomes Project)
  (7) ABT (Desmond Tutu)
  (8) KB1, Khoisan/Bushmen individual

# Counting SNPs

http://genome.ucsc.edu/index.html

# Counting SNPs

# Basic Shell Commands

```
$ cd [User_Folder]
$ mkdir session14
$ cd session14
```

# Counting SNPs

```
$ cp /home/biguser/tutor/session14/snp.txt .
$ less snp.txt
```

```
3263  A  A  A  A  A  A  A  A  T
3351  T  W  W  W  T  T  W  W  A
3544  T  T  T  T  T  T  T  Y  T
3567  T  T  T  T  T  T  T  Y  T
3774  K  G  T  T  T  T  T  T  T
4131  G  K  G  G  G  G  G  G  T
4190  A  A  A  R  A  A  A  R  A
4306  T  T  T  T  T  T  T  T  C
4371  C  Y  Y  Y  C  C  C  C  C
4489  G  R  A  A  A  A  A  A  A
6394  T  T  T  T  T  T  T  T  C
6523  G  R  A  R  A  A  R  A  A
```

# Counting SNPs

```
38357  T  T  Y  T  T  T  Y  Y  T
38368  G  G  G  G  G  G  G  G  C
38392  T  T  T  T  T  T  T  T  C
```

- The first column is position
- Nucleotides from eight individuals
- The last column is the nucleotide of chimpanzee
- Positions where at least one genome has an unknown base have been removed
- Positions containing the same nucleotide in all nine genomes have been removed

# Counting SNPs

```
3263  A  A  A  A  A  A  A  A  T
3351  T  W  W  W  T  T  W  W  A
3544  T  T  T  T  T  T  T  Y  T
3567  T  T  T  T  T  T  T  Y  T
3774  K  G  T  T  T  T  T  T  T
4131  G  K  G  G  G  G  G  G  T
4190  A  A  A  R  A  A  A  R  A
4306  T  T  T  T  T  T  T  C
```

| IUPAC nucleotide code | Base |
|---|---|
| A | Adenine |
| C | Cytosine |
| G | Guanine |
| T (or U) | Thymine (or Uracil) |
| R | A or G |
| Y | C or T |
| S | G or C |
| W | A or T |
| K | G or T |
| M | A or C |
| B | C or G or T |
| D | A or G or T |
| H | A or C or T |
| V | A or C or G |
| N | any base |
| . or - | gap |

# Counting SNPs

Create a distance matrix from SNPs of 9 genomes

$ vi snp.py

```
 3 # obtain pairwise distances from snp data,
 4 # counting sites where at least one allele is different
 5
 6 import re
 7
 8 humans = [
 9
10     # SNPs appear in the SNP data file in columns in this order
11
12     'YH',          # Han chinese
13     'SJK',         # Seong-Jin Kim
14     'JW',          # James Watson
15     'CV',          # Craig Venter
16     'NA18507',     # Yoruban of 1000 Genomes project
17     'NA12891',     # Of Central European origin
18     'ABT',         # Archbishop Desmond Tutu
19     'KB1',         # Bushmen individual
20     'chimp'        # chimpanzee
21 ]
```

# Counting SNPs

```
24 # 1 #
25 # initialize the distance matrix with zero values
26 # for the diagonal cells
27
28 diff = []
29
30 for i in range(0, 10):
31     diff.append([])
32     for j in range(0, 10):
33         diff[i].append(0)
34 print(diff)
```

```
$ python snp.py
```

```
[[0, 0, 0, 0, 0, 0, 0, 0, 0, 0], [0, 0, 0, 0, 0, 0, 0, 0, 0, 0],
 [0, 0, 0, 0, 0, 0, 0, 0, 0, 0], [0, 0, 0, 0, 0, 0, 0, 0, 0, 0],
 [0, 0, 0, 0, 0, 0, 0, 0, 0, 0], [0, 0, 0, 0, 0, 0, 0, 0, 0, 0],
 [0, 0, 0, 0, 0, 0, 0, 0, 0, 0], [0, 0, 0, 0, 0, 0, 0, 0, 0, 0],
 [0, 0, 0, 0, 0, 0, 0, 0, 0, 0], [0,_0, 0, 0, 0, 0, 0, 0, 0, 0]]
```

# Counting SNPs

3263 A A A A A A A T
      i     j

```python
1  # read the snp data from file
2
3  for line in open('snp.txt'):
4      line = line.rstrip()
5      columns = re.split(' ', line)
6
7      # 2 #
8      for i in range(1, 9):
9          for j in range(i + 1, 10):
10
11             # 3 #
12
13             if columns[i] != columns[j]:
14                 diff[i][j] += 1
15
16                 # 4 #
17                 # to produce a symmetric matrix
18                 diff[j][i] += 1
19
20 print(diff)
```

$ python snp.py

[[0, 0, 0, 0, 0, 0, 0, 0, 0, 0], [0, 0, 44597, 53594, 53913, 67914, 53710, 68837, 77272, 593367],
 [0, 44597, 0, 54192, 54537, 68826, 55281, 69404, 76929, 593496], [0, 53594, 54192, 0, 50859, 702
84, 51260, 70256, 77590, 592751], [0, 53913, 54537, 50859, 0, 70149, 51009, 69659, 77369, 592632]
, [0, 67914, 68826, 70284, 70149, 0, 69245, 70057, 79508, 599102], [0, 53710, 55281, 51260, 51009
, 69245, 0, 69941, 78130, 594831], [0, 68837, 69404, 70256, 69659, 70057, 69941, 0, 77707, 599292
], [0, 77272, 76929, 77590, 77369, 79508, 78130, 77707, 0, 600776], [0, 593367, 593496, 592751, 5
92632, 599102, 594831, 599292, 600776, 0]]

# Counting SNPs

```python
1  # 5 #
2  # print a header for PHYLIP format
3  # with the number of species
4
5  print('   ', '9')
6
7  # print the matrix data
8
9  for i in range(1, 10):
10
11     # 6 #
12
13     txt = humans[i - 1]
14     txt = txt[0:7]
15     print(txt, end = ' ')
16     length = 10 - len(txt)
17     short = ' ' * (length - 2)
18     print(short, end = ' ')
19     for j in range(1, 10):
20         print(diff[i][j], end = ' ')
21
22     print('')
```

# Counting SNPs

```
$ python snp.py
```

```
[biguser@R440 session14]$ python snp.py
    9
YH         0 44597 53594 53913 67914 53710 68837 77272 593367
SJK        44597 0 54192 54537 68826 55281 69404 76929 593496
JW         53594 54192 0 50859 70284 51260 70256 77590 592751
CV         53913 54537 50859 0 70149 51009 69659 77369 592632
NA18507    67914 68826 70284 70149 0 69245 70057 79508 599102
NA12891    53710 55281 51260 51009 69245 0 69941 78130 594831
ABT        68837 69404 70256 69659 70057 69941 0 77707 599292
KB1        77272 76929 77590 77369 79508 78130 77707 0 600776
chimp      593367 593496 592751 592632 599102 594831 599292 600776 0
```

```
$ python snp.py > snp.out
```

# Phylip package - neighbor

- Phylip package

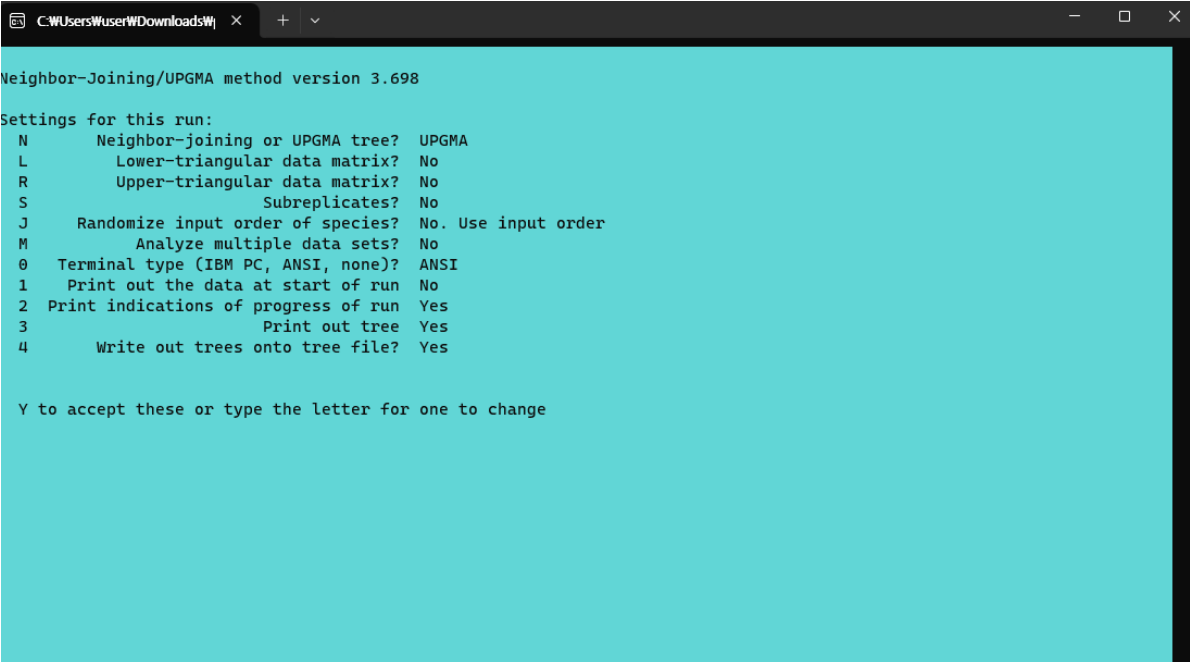  (https://phylipweb.github.io/phylip/)

# Phylip package - neighbor

```
neighbor: can't find input file "infile"
Please enter a new file name> snp.out
```

```
Neighbor-Joining/UPGMA method version 3.698

Settings for this run:
  N       Neighbor-joining or UPGMA tree?  UPGMA
  L          Lower-triangular data matrix?  No
  R          Upper-triangular data matrix?  No
  S                      Subreplicates?  No
  J    Randomize input order of species?  No. Use input order
  M          Analyze multiple data sets?  No
  0   Terminal type (IBM PC, ANSI, none)?  ANSI
  1    Print out the data at start of run  No
  2  Print indications of progress of run  Yes
  3                      Print out tree  Yes
  4       Write out trees onto tree file?  Yes


  Y to accept these or type the letter for one to change
```

# Phylip package - neighbor

```
Neighbor-Joining/UPGMA method version 3.698

Settings for this run:
  N         Neighbor-joining or UPGMA tree?  UPGMA
  L          Lower-triangular data matrix?  No
  R          Upper-triangular data matrix?  No
  S                    Subreplicates?  No
  J    Randomize input order of species?  No. Use input order
  M          Analyze multiple data sets?  No
  0    Terminal type (IBM PC, ANSI, none)?  ANSI
  1    Print out the data at start of run  No
  2    Print indications of progress of run  Yes
  3                    Print out tree  Yes
  4       Write out trees onto tree file?  Yes


 Y to accept these or type the letter for one to change
Y
```

```
Cycle    8: species 1 (22298.50000) joins species 2 (22298.50000)
Cycle    7: species 3 (25429.50000) joins species 4 (25429.50000)
Cycle    6: node 3 ( 137.75000) joins species 6 (25567.25000)
Cycle    5: node 1 (4803.75000) joins node 3 (1535.00000)
Cycle    4: node 1 (7539.55000) joins species 5 (34641.80000)
Cycle    3: node 1 ( 204.36667) joins species 7 (34846.16667)
Cycle    2: node 1 (4047.04762) joins species 8 (38893.21429)
Cycle    1: node 1 (258997.22321) joins species 9 (297890.43750)

Output written on file "outfile"

Tree written on file "outtree"

Done.
```

## outfile



```
 9 Populations

Neighbor-Joining/UPGMA method version 3.698


UPGMA method

Negative branch lengths allowed

                                              +---YH
                                            +-1
                                            ! +---SJK
                                          +-4
                                          ! !    +----JW
                                          ! ! +-2
                                        +-5 +-3 +----CV
                                        ! ! !
                                        ! !  +----NA12891
                                      +-6 !
                                      ! ! +------NA18507
+-------------------------------------7 !
!                                     ! +------ABT
--8                                   !
 !                                    +-------KB1
 !
 +-------------------------------------------------------chimp

From    To        Length        Height
----    --        ------        ------
  8      7      258997.22321   258997.22321
  7      6      4047.04762     263044.27083
  6      5       204.36667     263248.63750
  5      4      7539.55000     270788.18750
  4      1      4803.75000     275591.93750
  1     YH      22298.50000    297890.43750
  1     SJK     22298.50000    297890.43750
  4      3      1535.00000     272323.18750
  3      2       137.75000     272460.93750
  2     JW      25429.50000    297890.43750
  2     CV      25429.50000    297890.43750
  3    NA12891  25567.25000    297890.43750
  5    NA18507  34641.80000    297890.43750
  6     ABT     34846.16667    297890.43750
  7     KB1     38893.21429    297890.43750
  8    chimp    297890.43750   297890.43750
```

## outtree

```
((((((YH:22298.50000,SJK:22298.50000):4803.75000,((JW:25429.50000,
CV:25429.50000):137.75000,NA12891:25567.25000):1535.00000):7539.55000,
NA18507:34641.80000):204.36667,ABT:34846.16667):4047.04762,
KB1:38893.21429):258997.22321,chimp:297890.43750);
```

# NJplot

- NJplot

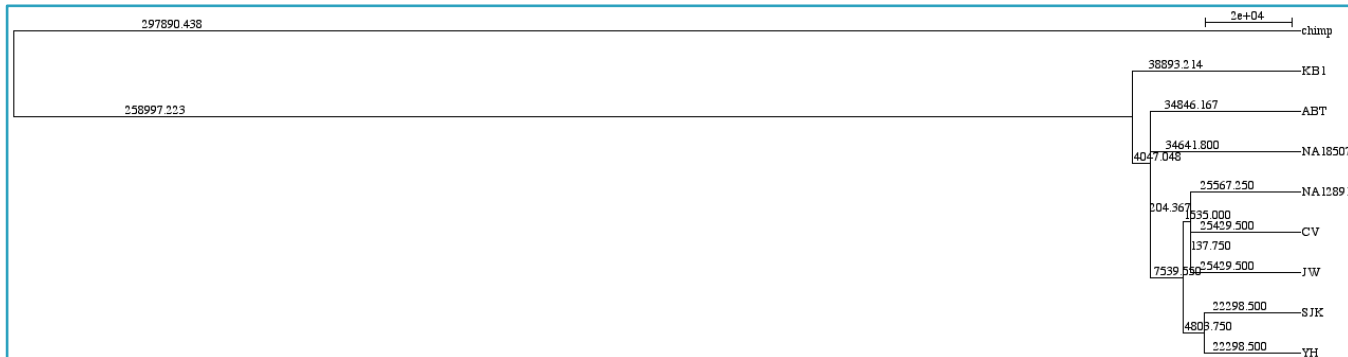  (http://doua.prabi.fr/software/njplot)

  **NJplot**

  NEW: NJplot plots trees in PDF and PostScript formats (not for MacOS).
  NEW: NJplot allows to open several tree windows.
  NEW: NJplot can draw multibranched trees with or without branch lengths.

  NJplot is a tree drawing program able to draw any phylogenetic tree expressed in the Newick phylogenetic tree format (*e.g.*, the format used by the PHYLIP package). NJplot is especially convenient for rooting the unrooted trees obtained from parsimony, distance or maximum likelihood tree-building methods.
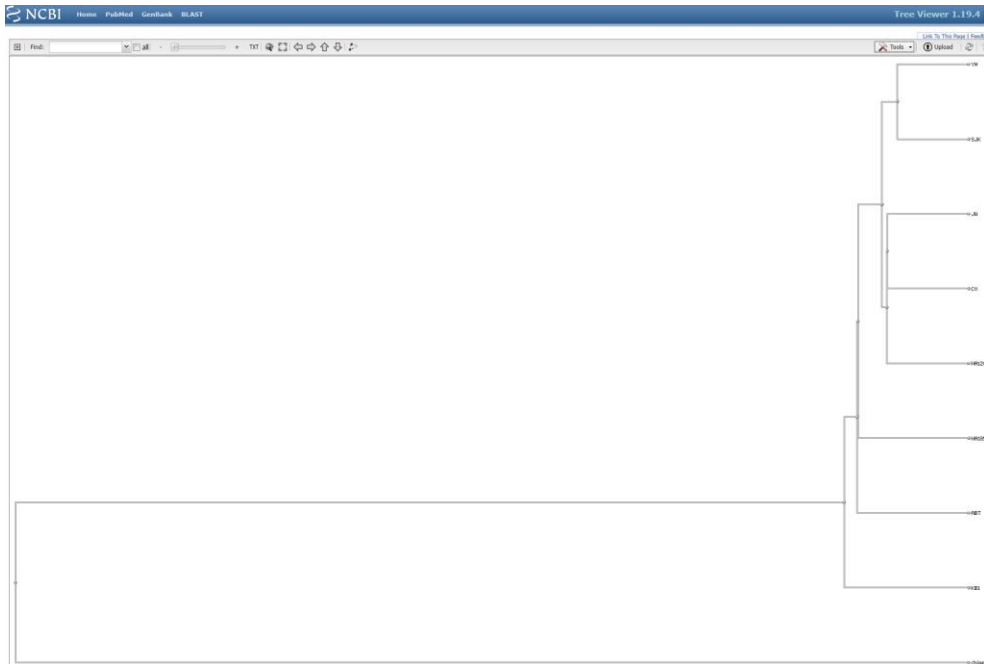
  A screen shot of the main window of njplot is available here.

# NCBI tree viewer

- NCBI tree viewer

  (https://www.ncbi.nlm.nih.gov/projects/treeview/)

# Exercise

□ The file 'chr4snp.txt' is a list of SNPs in the human chromosome 4, according to dbSNP build 130. Write a Python script that will list the SNPs (positions) that are present in this file but that are not found in the file 'snp.txt'. The file 'chr4snp.txt' uses 'zero-based' numbering. From a practical point of view, this means that the third column positions in that file are comparable to the position numbers in 'snp.txt'.

```
cp /home/biguser/tutor/session14/chr4snp.txt .
```

# Exercise



Snp.txt

실제 position  chr4snp.txt

Chr4snp.txt 에 있는 snp 들중 snp.txt에 없는 position을 출력 할 것!

# Exercise

```python
1 # exercise
2
3 import sys
4
5 infile1 = open(sys.argv[1], 'r') # snp.txt
6 infile2 = open(sys.argv[2], 'r') # chr4snp.txt
7
8 snppos = dict()
9 for line in infile1.readlines():
10        col = line.split(' ')
11        pos = col[0]
12        snppos[pos] = ''
13 infile1.close()
14
15 for line in infile2.readlines():
16        line = line.strip()
17        if not line.startswith('#'):
18            col = line.split('\t')
19            pos = col[2]
20            if not snppos.has_key(pos):
21                print(pos)
22 infile2.close()
```