# class10

Mina Wu( PID: A59013200)

2024-02-09

## What is in the PDB?

Download a CSV file with current composiiton data from: https://www.rcsb.org/stats/summary

```
pdbstats <- read.csv("Data Export Summary.csv", row.names=1)
pdbstats
```

|  | X.ray | EM | NMR | Multiple.methods | Neutron | Other |
|---|---|---|---|---|---|---|
| Protein (only) | 161,663 | 12,592 | 12,337 | 200 | 74 | 32 |
| Protein/Oligosaccharide | 9,348 | 2,167 | 34 | 8 | 2 | 0 |
| Protein/NA | 8,404 | 3,924 | 286 | 7 | 0 | 0 |
| Nucleic acid (only) | 2,758 | 125 | 1,477 | 14 | 3 | 1 |
| Other | 164 | 9 | 33 | 0 | 0 | 0 |
| Oligosaccharide (only) | 11 | 0 | 6 | 1 | 0 | 4 |

|  | Total |
|---|---|
| Protein (only) | 186,898 |
| Protein/Oligosaccharide | 11,559 |
| Protein/NA | 12,621 |
| Nucleic acid (only) | 4,378 |
| Other | 206 |
| Oligosaccharide (only) | 22 |

```
#x<- "2,2222"

#sub <- as.numeric(gsub(",", "", x))
#sub
```

Create a function to remove commas.

```
commasum<- function (x){
  #Remove comma, convert to numeric and sum
 sum(as.numeric(gsub(",", "", x)))
}
```

```
commasum(pdbstats$X.ray)
```

[1] 182348

I can now 'apply()' this function to my wee table to get all the numbers I need.

```
round(apply(pdbstats,2, commasum)/commasum(pdbstats$Total)*100,2)
```

```
        X.ray              EM             NMR Multiple.methods
        84.54            8.72            6.57             0.11
      Neutron           Other           Total
        0.04            0.02          100.00
```

Q1: What percentage of structures in the PDB are solved by X-Ray and Electron Microscopy.

84.54% for xray and 8.72% for EM.

```
commasum(pdbstats$Total[1])/commasum(pdbstats$Total)
```

[1] 0.8665362

Q2: What proportion of structures in the PDB are protein?

0.8665362

Q3 class: How does the total number of protein structures in the PDB relate to the total number of protein sequences in Uniprot?

(186898/250,322,721)*100 = 0.074%

Q3: Type HIV in the PDB website search box on the home page and determine how many HIV-1 protease structures are in the current PDB?

There ar 4,410 structures.

# Visualizing the HIV-1 potease struture

We will use the Mol* (mol-star) homepage at: https://molstar.org/viewer/
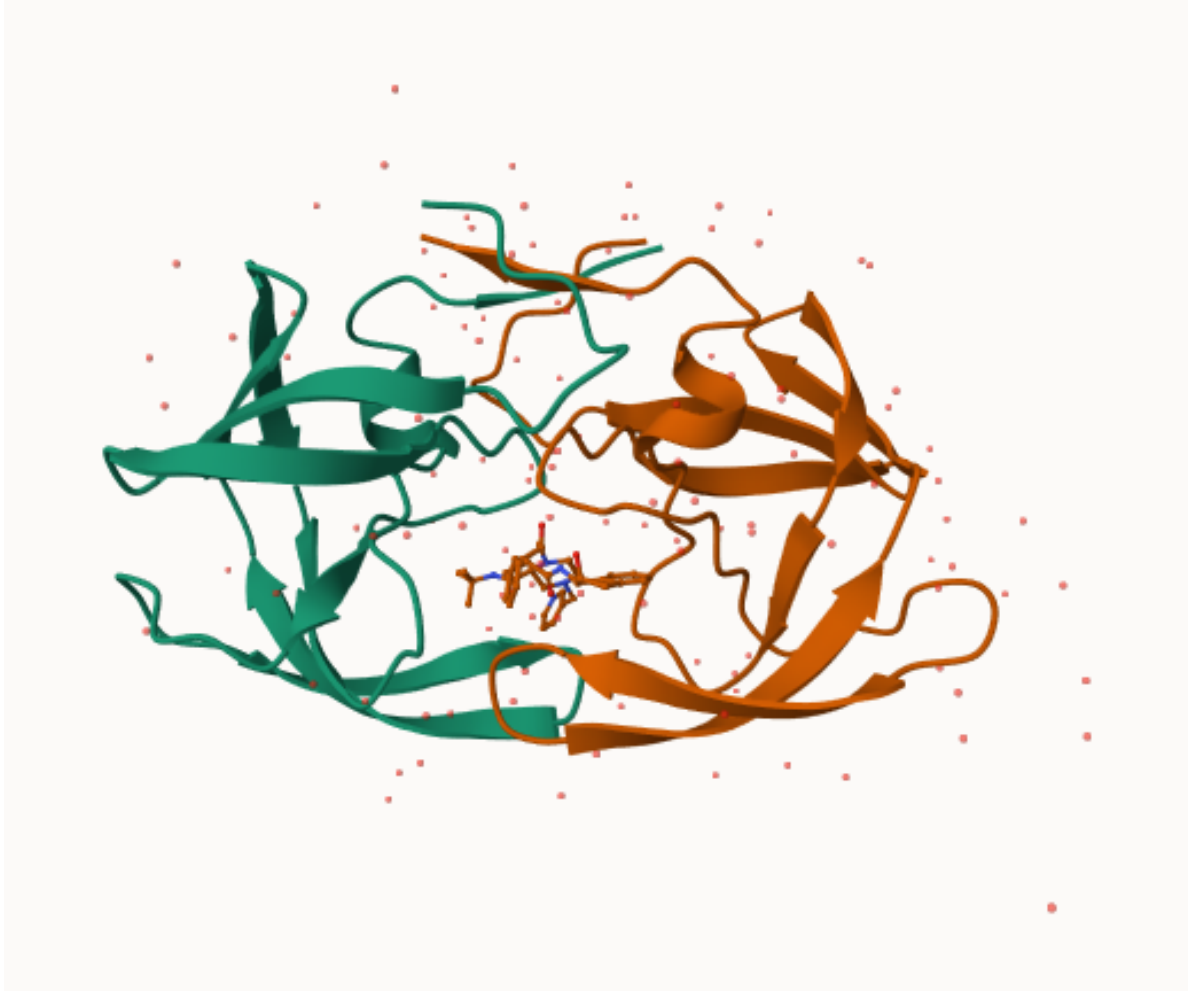
Looking at 1HSG



Figure 1: Figure 1HSG

Looking at the aspartic acid at position 25 for both the A and B chain. It's represented as space-fill.

## Working with structures in R

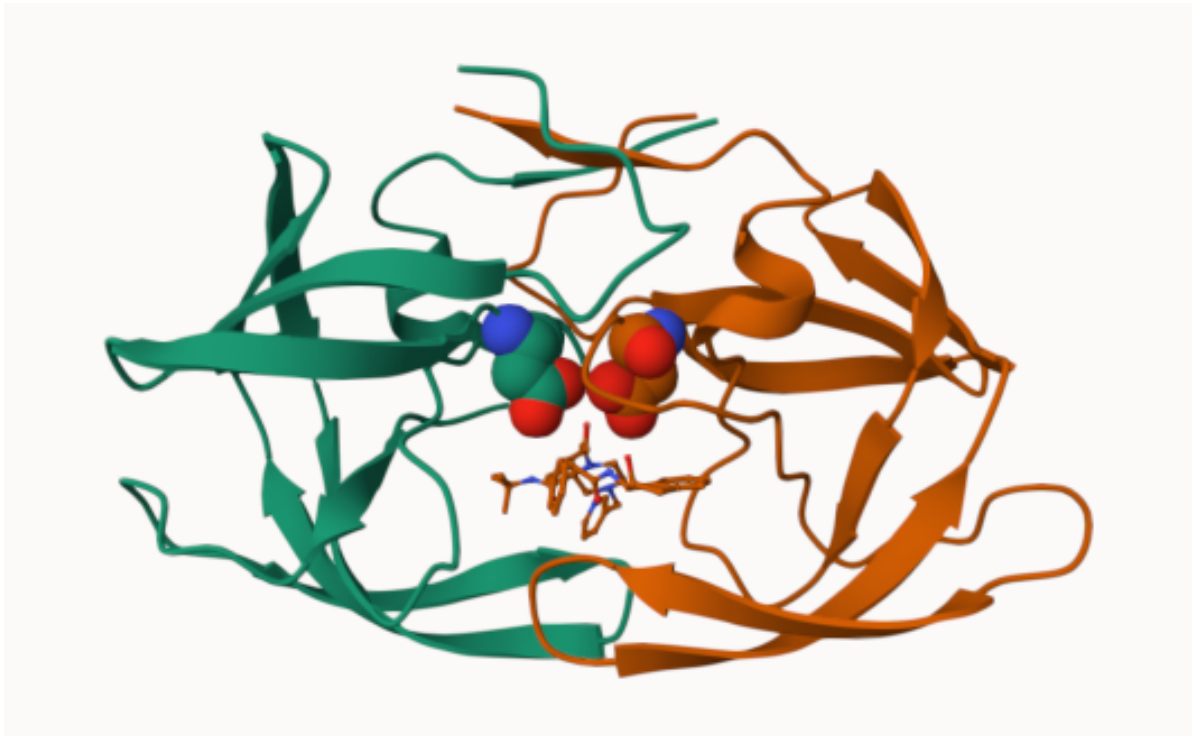We will use the bio3d package for structrual bioinformatics

Figure 2: Figure 1HSG

```r
library(bio3d)
```

Warning: package 'bio3d' was built under R version 4.0.5

```r
hiv<- read.pdb("1hsg")
```

Note: Accessing on-line PDB file

```r
hiv
```

Call:  read.pdb(file = "1hsg")

   Total Models#: 1
     Total Atoms#: 1686,  XYZs#: 5058  Chains#: 2  (values: A B)

     Protein Atoms#: 1514  (residues/Calpha atoms#: 198)
     Nucleic acid Atoms#: 0  (residues/phosphate atoms#: 0)

     Non-protein/nucleic Atoms#: 172  (residues: 128)
     Non-protein/nucleic resid values: [ HOH (127), MK1 (1) ]

   Protein sequence:
      PQITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWKPKMIGGIGGFIKVRQYD
      QILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNFPQITLWQRPLVTIKIGGQLKE
      ALLDTGADDTVLEEMSLPGRWKPKMIGGIGGFIKVRQYDQILIEICGHKAIGTVLVGPTP
      VNIIGRNLLTQIGCTLNF

+ attr: atom, xyz, seqres, helix, sheet,
        calpha, remark, call

```r
head (hiv$atom)
```

```
  type eleno elety  alt resid chain resno insert      x      y     z o     b
1 ATOM     1     N <NA>   PRO     A     1   <NA> 29.361 39.686 5.862 1 38.10
2 ATOM     2    CA <NA>   PRO     A     1   <NA> 30.307 38.663 5.319 1 40.62
3 ATOM     3     C <NA>   PRO     A     1   <NA> 29.760 38.071 4.022 1 42.64
```

```
4 ATOM      4      O <NA>    PRO     A     1    <NA> 28.600 38.302 3.676 1 43.40
5 ATOM      5     CB <NA>    PRO     A     1    <NA> 30.508 37.541 6.342 1 37.87
6 ATOM      6     CG <NA>    PRO     A     1    <NA> 29.296 37.591 7.162 1 38.40
  segid elesy charge
1  <NA>     N   <NA>
2  <NA>     C   <NA>
3  <NA>     C   <NA>
4  <NA>     O   <NA>
5  <NA>     C   <NA>
6  <NA>     C   <NA>
```

```
  pdbseq(hiv)
```

```
  1    2    3    4    5    6    7    8    9   10   11   12   13   14   15   16   17   18   19   20
"P"  "Q"  "I"  "T"  "L"  "W"  "Q"  "R"  "P"  "L"  "V"  "T"  "I"  "K"  "I"  "G"  "G"  "Q"  "L"  "K"
 21   22   23   24   25   26   27   28   29   30   31   32   33   34   35   36   37   38   39   40
"E"  "A"  "L"  "L"  "D"  "T"  "G"  "A"  "D"  "D"  "T"  "V"  "L"  "E"  "E"  "M"  "S"  "L"  "P"  "G"
 41   42   43   44   45   46   47   48   49   50   51   52   53   54   55   56   57   58   59   60
"R"  "W"  "K"  "P"  "K"  "M"  "I"  "G"  "G"  "I"  "G"  "G"  "F"  "I"  "K"  "V"  "R"  "Q"  "Y"  "D"
 61   62   63   64   65   66   67   68   69   70   71   72   73   74   75   76   77   78   79   80
"Q"  "I"  "L"  "I"  "E"  "I"  "C"  "G"  "H"  "K"  "A"  "I"  "G"  "T"  "V"  "L"  "V"  "G"  "P"  "T"
 81   82   83   84   85   86   87   88   89   90   91   92   93   94   95   96   97   98   99    1
"P"  "V"  "N"  "I"  "I"  "G"  "R"  "N"  "L"  "L"  "T"  "Q"  "I"  "G"  "C"  "T"  "L"  "N"  "F"  "P"
  2    3    4    5    6    7    8    9   10   11   12   13   14   15   16   17   18   19   20   21
"Q"  "I"  "T"  "L"  "W"  "Q"  "R"  "P"  "L"  "V"  "T"  "I"  "K"  "I"  "G"  "G"  "Q"  "L"  "K"  "E"
 22   23   24   25   26   27   28   29   30   31   32   33   34   35   36   37   38   39   40   41
"A"  "L"  "L"  "D"  "T"  "G"  "A"  "D"  "D"  "T"  "V"  "L"  "E"  "E"  "M"  "S"  "L"  "P"  "G"  "R"
 42   43   44   45   46   47   48   49   50   51   52   53   54   55   56   57   58   59   60   61
"W"  "K"  "P"  "K"  "M"  "I"  "G"  "G"  "I"  "G"  "G"  "F"  "I"  "K"  "V"  "R"  "Q"  "Y"  "D"  "Q"
 62   63   64   65   66   67   68   69   70   71   72   73   74   75   76   77   78   79   80   81
"I"  "L"  "I"  "E"  "I"  "C"  "G"  "H"  "K"  "A"  "I"  "G"  "T"  "V"  "L"  "V"  "G"  "P"  "T"  "P"
 82   83   84   85   86   87   88   89   90   91   92   93   94   95   96   97   98   99
"V"  "N"  "I"  "I"  "G"  "R"  "N"  "L"  "L"  "T"  "Q"  "I"  "G"  "C"  "T"  "L"  "N"  "F"
```

```
  aa123(pdbseq(hiv)[25])
```

```
[1] "ASP"
```

```
  adk <- read.pdb("6s36")
```

```
 Note: Accessing on-line PDB file
  PDB has ALT records, taking A only, rm.alt=TRUE
```

  adk

```
 Call:  read.pdb(file = "6s36")

   Total Models#: 1
     Total Atoms#: 1898,  XYZs#: 5694  Chains#: 1  (values: A)

     Protein Atoms#: 1654  (residues/Calpha atoms#: 214)
     Nucleic acid Atoms#: 0  (residues/phosphate atoms#: 0)

     Non-protein/nucleic Atoms#: 244  (residues: 244)
     Non-protein/nucleic resid values: [ CL (3), HOH (238), MG (2), NA (1) ]

   Protein sequence:
      MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMLRAAVKSGSELGKQAKDIMDAGKLVT
      DELVIALVKERIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFDVPDELIVDKI
      VGRRVHAPSGRVYHVKFNPPKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQMTAPLIG
      YYSKEAEAGNTKYAKVDGTKPVAEVRADLEKILG

+ attr: atom, xyz, seqres, helix, sheet,
        calpha, remark, call
```
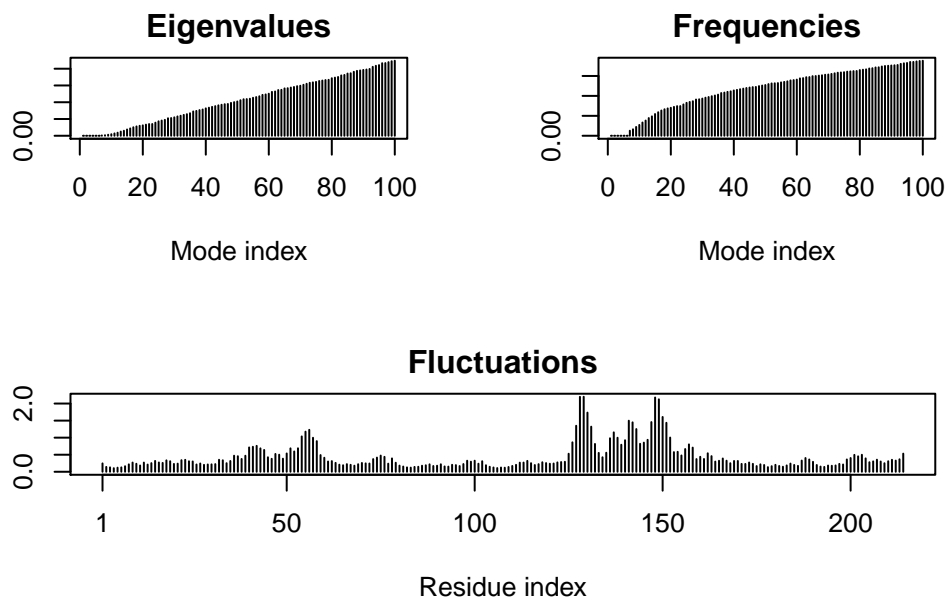
NOrmal mode analysis (NMA) a bioinformatics method to predict fucntional motions and large-scale structure changes.

  m <- nma(adk)

```
 Building Hessian...        Done in 0.042 seconds.
 Diagonalizing Hessian...   Done in 0.503 seconds.
```

  plot(m)

Make a wee move (a.k.a "trajectory") of this predicted motions

```
mktrj(m, file="adk_movie.pdb")
```

## Quick comparactive Structure analysis of Adenylate kinase

Extract sequence and run a BLAST search

```
s<- pdbseq(adk)
blast<- blast.pdb(s)
```

```
Searching ... please wait (updates every 5 seconds) RID = WDH5K2Y6016
.....
Reporting 83 hits
```
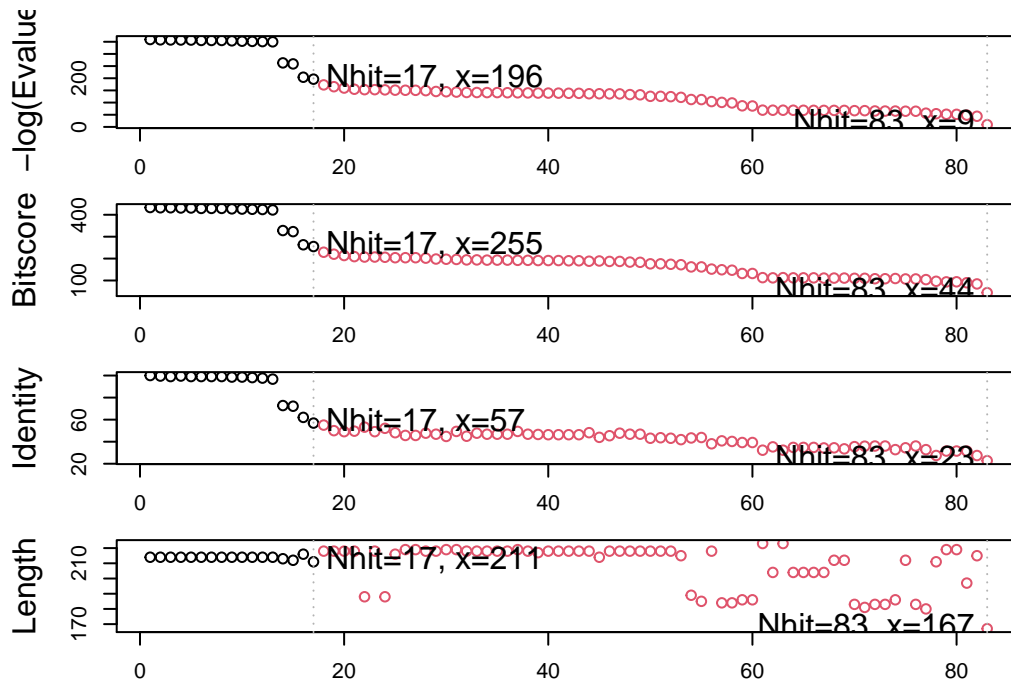
```
hits<- plot(blast)
```

```
* Possible cutoff values:    196 9
          Yielding Nhits:    17 83
```

```
* Chosen cutoff value of:      196
          Yielding Nhits:      17
```



Get the results from BLAST and download all the top hits.

```
hits$pdb.id
```

```
 [1] "6S36_A" "1AKE_A" "8BQF_A" "6RZE_A" "4X8M_A" "4X8H_A" "1E4V_A" "3HPR_A"
 [9] "5EJE_A" "1E4Y_A" "3X2S_A" "6HAP_A" "6HAM_A" "4K46_A" "4NP6_A" "3GMT_A"
[17] "4PZL_A"
```

```
# Download releated PDB files
files <- get.pdb(hits$pdb.id, path="pdbs", split=TRUE, gzip=TRUE)
```

```
Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/6S36.pdb.gz exists. Skipping download
```

```
Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/1AKE.pdb.gz exists. Skipping download
```

9

```
Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/8BQF.pdb.gz exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/6RZE.pdb.gz exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/4X8M.pdb.gz exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/4X8H.pdb.gz exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/1E4V.pdb.gz exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/3HPR.pdb.gz exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/5EJE.pdb.gz exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/1E4Y.pdb.gz exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/3X2S.pdb.gz exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/6HAP.pdb.gz exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/6HAM.pdb.gz exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/4K46.pdb.gz exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/4NP6.pdb.gz exists. Skipping download
```

```
Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/3GMT.pdb.gz exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/4PZL.pdb.gz exists. Skipping download


  |
  |                                                                  |   0%
  |
  |====                                                              |   6%
  |
  |=======                                                           |  12%
  |
  |===========                                                       |  18%
  |
  |===============                                                   |  24%
  |
  |===================                                               |  29%
  |
  |=======================                                           |  35%
  |
  |===========================                                       |  41%
  |
  |==============================                                    |  47%
  |
  |==================================                                |  53%
  |
  |======================================                            |  59%
  |
  |==========================================                        |  65%
  |
  |==============================================                    |  71%
  |
  |==================================================                |  76%
  |
  |======================================================            |  82%
  |
  |==========================================================        |  88%
  |
  |==============================================================    |  94%
  |
```

```
  |========================================================================| 100%
```

```
  # Align releated PDBs
  pdbs <- pdbaln(files, fit = TRUE, exefile="msa")
```

```
Reading PDB files:
pdbs/split_chain/6S36_A.pdb
pdbs/split_chain/1AKE_A.pdb
pdbs/split_chain/8BQF_A.pdb
pdbs/split_chain/6RZE_A.pdb
pdbs/split_chain/4X8M_A.pdb
pdbs/split_chain/4X8H_A.pdb
pdbs/split_chain/1E4V_A.pdb
pdbs/split_chain/3HPR_A.pdb
pdbs/split_chain/5EJE_A.pdb
pdbs/split_chain/1E4Y_A.pdb
pdbs/split_chain/3X2S_A.pdb
pdbs/split_chain/6HAP_A.pdb
pdbs/split_chain/6HAM_A.pdb
pdbs/split_chain/4K46_A.pdb
pdbs/split_chain/4NP6_A.pdb
pdbs/split_chain/3GMT_A.pdb
pdbs/split_chain/4PZL_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
.   PDB has ALT records, taking A only, rm.alt=TRUE
.   PDB has ALT records, taking A only, rm.alt=TRUE
.   PDB has ALT records, taking A only, rm.alt=TRUE
....    PDB has ALT records, taking A only, rm.alt=TRUE
.   PDB has ALT records, taking A only, rm.alt=TRUE
....    PDB has ALT records, taking A only, rm.alt=TRUE
.   PDB has ALT records, taking A only, rm.alt=TRUE
....

Extracting sequences

pdb/seq: 1   name: pdbs/split_chain/6S36_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 2   name: pdbs/split_chain/1AKE_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 3   name: pdbs/split_chain/8BQF_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
```
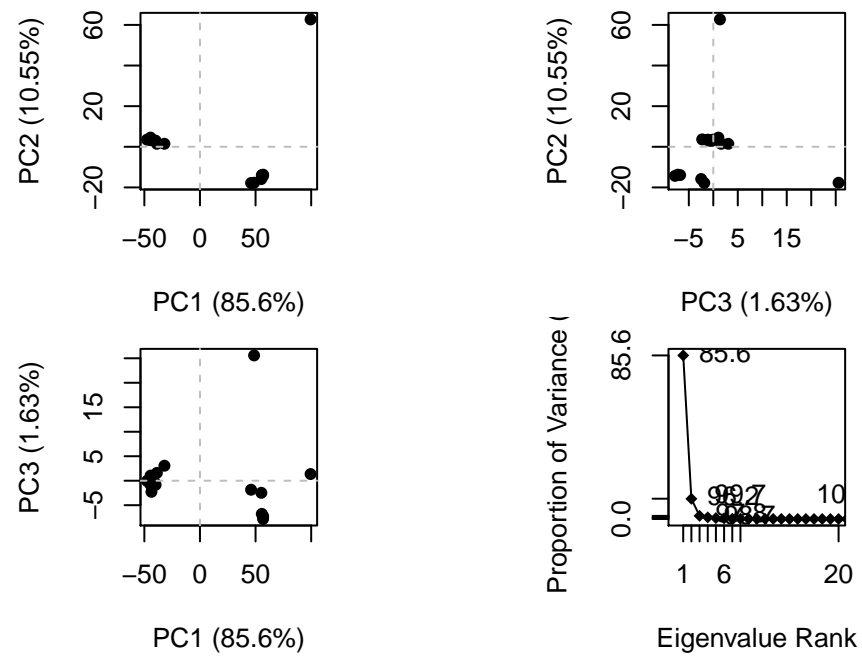
```
pdb/seq: 4    name: pdbs/split_chain/6RZE_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 5    name: pdbs/split_chain/4X8M_A.pdb
pdb/seq: 6    name: pdbs/split_chain/4X8H_A.pdb
pdb/seq: 7    name: pdbs/split_chain/1E4V_A.pdb
pdb/seq: 8    name: pdbs/split_chain/3HPR_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 9    name: pdbs/split_chain/5EJE_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 10   name: pdbs/split_chain/1E4Y_A.pdb
pdb/seq: 11   name: pdbs/split_chain/3X2S_A.pdb
pdb/seq: 12   name: pdbs/split_chain/6HAP_A.pdb
pdb/seq: 13   name: pdbs/split_chain/6HAM_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 14   name: pdbs/split_chain/4K46_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 15   name: pdbs/split_chain/4NP6_A.pdb
pdb/seq: 16   name: pdbs/split_chain/3GMT_A.pdb
pdb/seq: 17   name: pdbs/split_chain/4PZL_A.pdb
```

## PCA of all these ADK structures

```
pc.xray <- pca(pdbs)
plot(pc.xray)
```

```
mktrj(pc.xray, file="pco_movie.pdb")
```