# Class10 : Structural Bioinformatics (Pt. 1)

Safiya Sayd (PID: A18027139)

2025-02-06

## Table of contents

## 1. PDB database

The main repository of biomolecular structior data is called the pdb found at: https://www.rcsb.org/

Let's see what this database contains. PDB analyze >PDB statistics > by Experimental Method and Molecular Type

```
pdbstats <- read.csv("Data Export Summary.csv")
pdbstats
```

|   | Molecular.Type | X.ray | EM | NMR | Multiple.methods | Neutron | Other |
|---|---|---|---|---|---|---|---|
| 1 | Protein (only) | 169,563 | 16,774 | 12,578 | 208 | 81 | 32 |
| 2 | Protein/Oligosaccharide | 9,939 | 2,839 | 34 | 8 | 2 | 0 |
| 3 | Protein/NA | 8,801 | 5,062 | 286 | 7 | 0 | 0 |
| 4 | Nucleic acid (only) | 2,890 | 151 | 1,521 | 14 | 3 | 1 |
| 5 | Other | 170 | 10 | 33 | 0 | 0 | 0 |
| 6 | Oligosaccharide (only) | 11 | 0 | 6 | 1 | 0 | 4 |

|   | Total |
|---|---|
| 1 | 199,236 |
| 2 | 12,822 |
| 3 | 14,156 |
| 4 | 4,580 |

1

```
5      213
6       22
```

commas result in numerical values being categorized as a character

```
pdbstats$X.ray
```

```
[1] "169,563" "9,939"   "8,801"   "2,890"   "170"      "11"
```

This can be fixed by replacing "," for nothing " " with **sub()** function:

```
x <- pdbstats$X.ray
x_numeric <- as.numeric( gsub(",", "", x))
x_numeric
```

```
[1] 169563   9939   8801   2890    170     11
```

or I can use the **readr** package and the **read_csv()** function

```
library(readr)
pdbstats <- read_csv("Data Export Summary.csv", show_col_types = FALSE)
pdbstats
```

```
# A tibble: 6 x 8
  `Molecular Type`   `X-ray`    EM    NMR `Multiple methods` Neutron Other   Total
  <chr>               <dbl> <dbl> <dbl>              <dbl>   <dbl> <dbl>   <dbl>
1 Protein (only)     169563 16774 12578                208      81    32 199236
2 Protein/Oligosacc~   9939  2839    34                  8       2     0  12822
3 Protein/NA           8801  5062   286                  7       0     0  14156
4 Nucleic acid (onl~   2890   151  1521                 14       3     1   4580
5 Other                 170    10    33                  0       0     0    213
6 Oligosaccharide (~     11     0     6                  1       0     4     22
```

I want to clean the column names so that they are all lower case and don't have spaces in them

```
colnames(pdbstats)
```

```
[1] "Molecular Type"    "X-ray"             "EM"              "NMR"
[5] "Multiple methods"  "Neutron"           "Other"           "Total"
```

```
library(janitor)
```

```
Attaching package: 'janitor'

The following objects are masked from 'package:stats':

    chisq.test, fisher.test
```

```
df <- clean_names(pdbstats)
df
```

```
# A tibble: 6 x 8
  molecular_type      x_ray    em   nmr multiple_methods neutron other  total
  <chr>               <dbl> <dbl> <dbl>            <dbl>   <dbl> <dbl>  <dbl>
1 Protein (only)     169563 16774 12578              208      81    32 199236
2 Protein/Oligosacchar~ 9939  2839    34                8       2     0  12822
3 Protein/NA           8801  5062   286                7       0     0  14156
4 Nucleic acid (only)  2890   151  1521               14       3     1   4580
5 Other                 170    10    33                0       0     0    213
6 Oligosaccharide (onl~  11     0     6                1       0     4     22
```

```
sum(df$x_ray)
```

```
[1] 191374
```

```
df$total
```

```
[1] 199236  12822  14156   4580    213     22
```

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

93.58566%

Percent of X-ray structures

```
sum(df$x_ray)/sum(df$total) *100
```

```
[1] 82.83549
```

Percent of EM structure

```
sum(df$em)/sum(df$total) *100
```

```
[1] 10.75017
```

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

86.23852

```
protein_row <- df[df$molecular_type =="Protein (only)",]
protein_row
```

```
# A tibble: 1 x 8
  molecular_type  x_ray     em    nmr multiple_methods neutron other   total
  <chr>           <dbl>  <dbl>  <dbl>            <dbl>   <dbl> <dbl>   <dbl>
1 Protein (only) 169563  16774  12578              208      81    32  199236
```

```
sum(protein_row$total) / sum(df$total) * 100
```

```
[1] 86.23852
```

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?

2,298 structures are currently in the pdb.

## 2. Using mol*

The main mol* homepage: https://molstar.org/viewer/ We can input our own PDB files or just give it a PDB database accession code (4 letter PDB code)
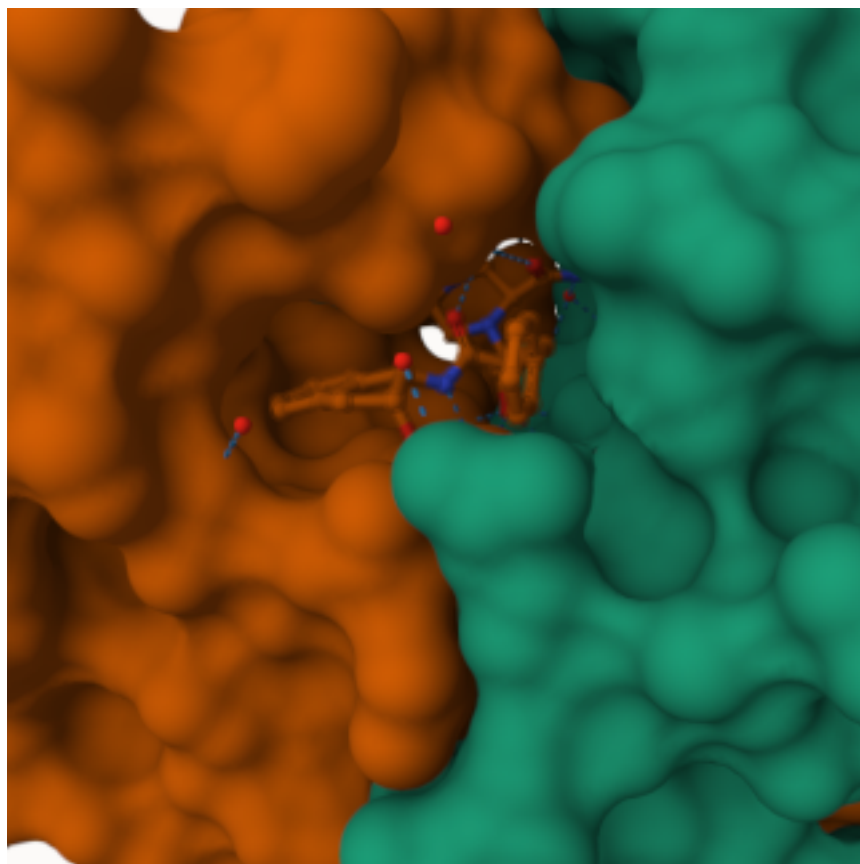
Figure 1: Molecular View of 1HSG

Q4: Water molecules normally have 3 atoms. Why do we see just one atom per water molecule in this structure?

This is due to the limitation of detection of water molecules. Hydrogen molecules are harder to detect while oxygen molecules aren't. The one atom shown is oxygen and has a higher accuracy of position. Another reason this is beneficial is it simplifies visualizing the structure.

Q5. There is a critical "conserved" water molecule in the binding site. Can you identify this water molecule? What residue number does this water molecule have
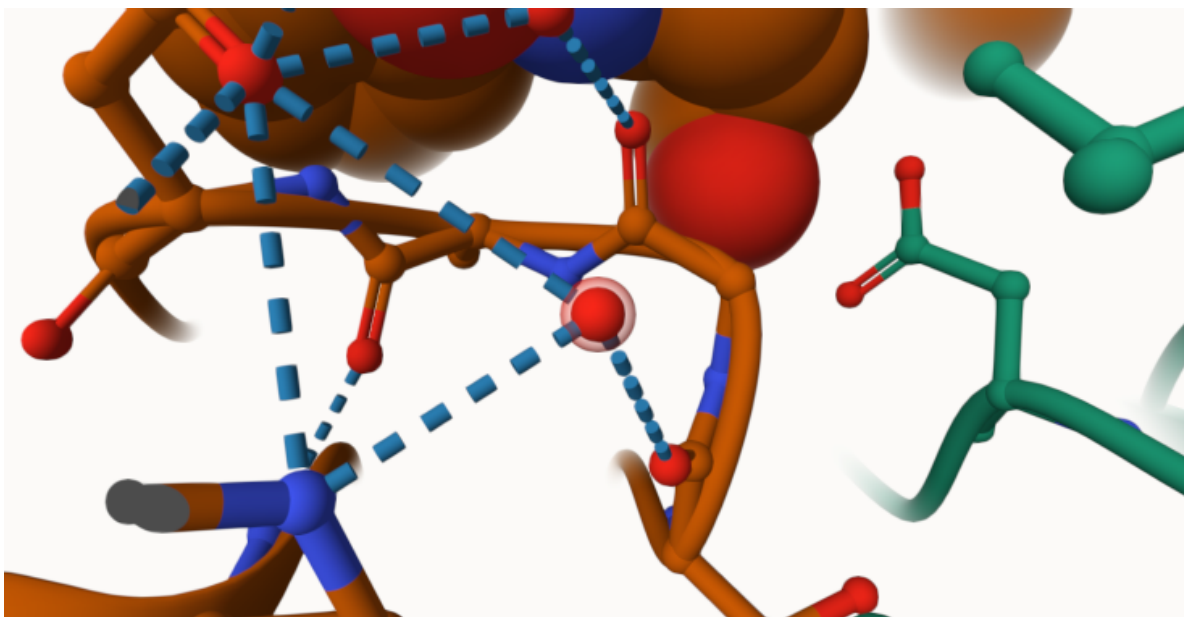
306

Figure 2: Conserved Water Molecule

Q6. Generate and save a figure clearly showing the two distinct chains of HIV-protease along with the ligand.
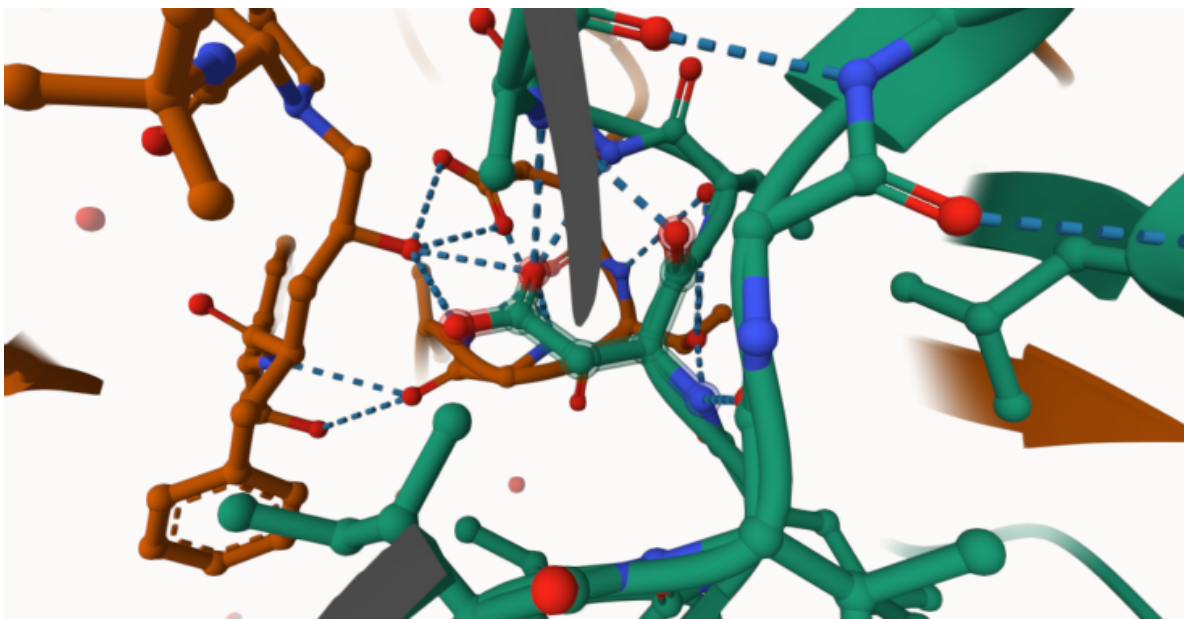


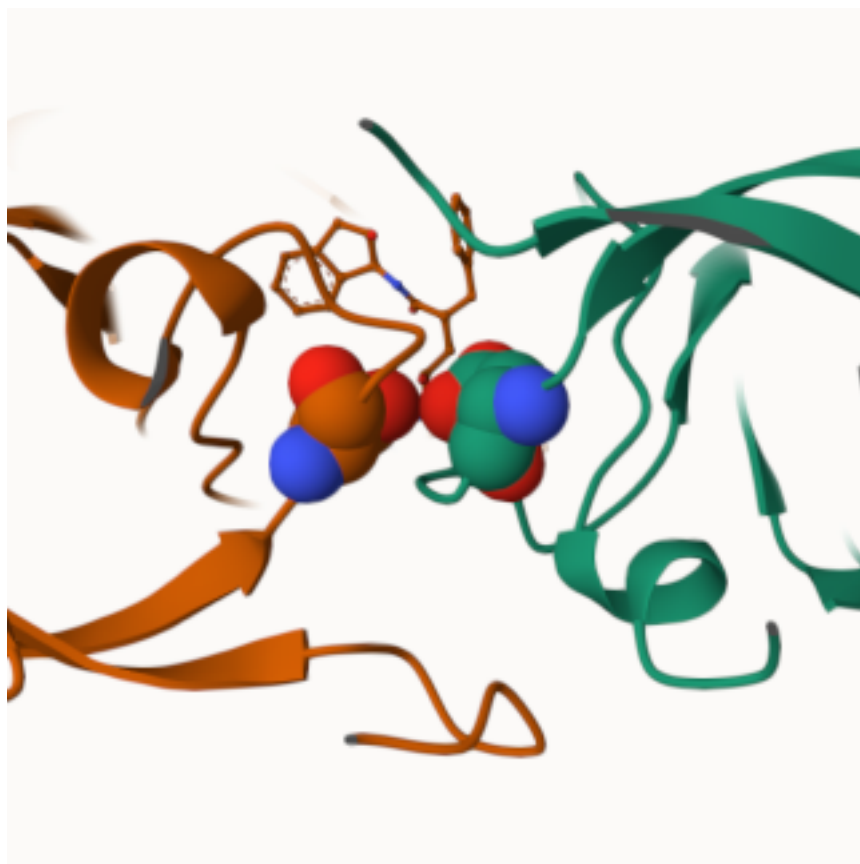Figure 3: aspartic acid interaction with ligand
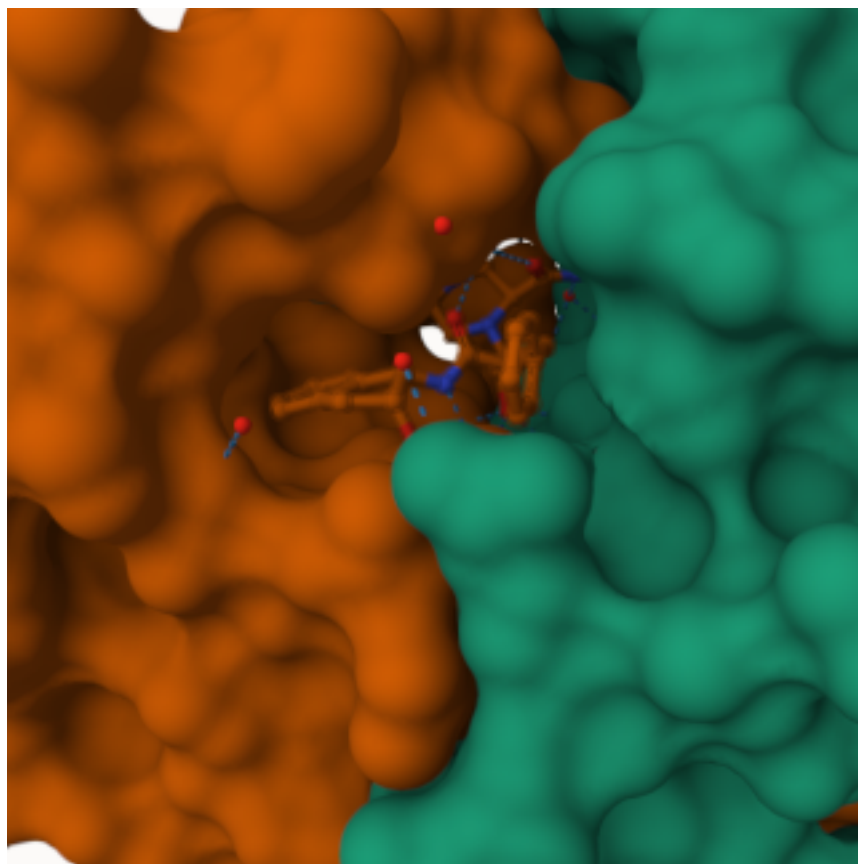
Figure 4: Aspartic Acid 25

Figure 5: Molecular View of 1HSG

## 3. Introduction to Bio3D in R

**Bio3D** can be used to read PDB data in R

```
library(bio3d)

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

```
  Note: Accessing on-line PDB file
```

```
pdb
```

```
 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
```

Q7. How many amino acid residues are there in this pdb object?

```r
length( pdbseq(pdb))
```

```
[1] 198
```

Q8. Name one of the two non-protein residues?

MK1

9. How many protein chains are in this structure?

2 chains, A and B

Looking in more detail of `pdb`

```r
attributes(pdb)
```

```
$names
[1] "atom"   "xyz"    "seqres" "helix"  "sheet"  "calpha" "remark" "call"

$class
[1] "pdb" "sse"
```

```
head(pdb$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>
```
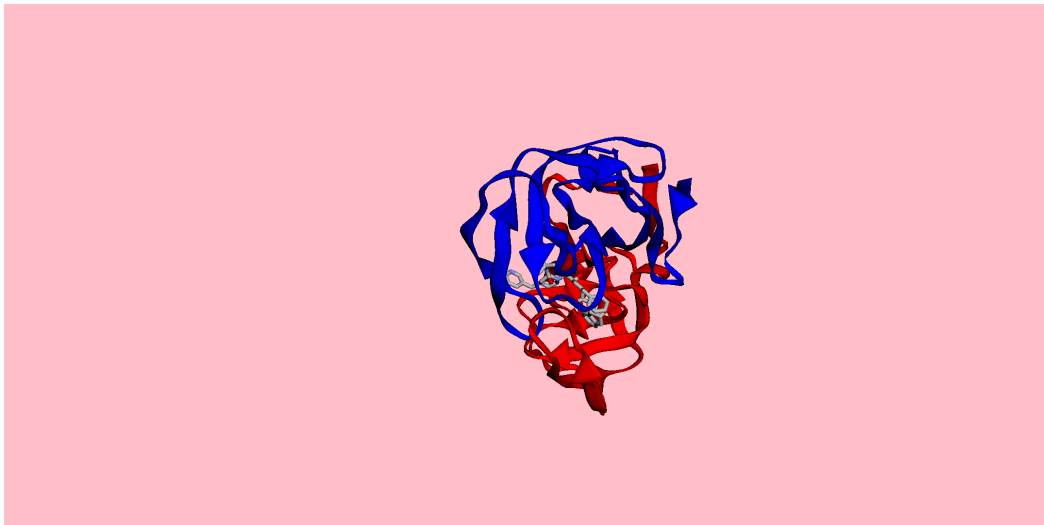
new function in Bio3D install.packages("r3dmol") and install.packages("shiny")

```
source("https://tinyurl.com/viewpdb")
view.pdb(pdb, backgroundColor = "pink")
```

```
file:////private/var/folders/sy/_fr9v5r51nxc7bzp4h_683nh0000gn/T/Rtmp5g03vi/filea38e75d7f61c/
```

10

## 4. predicting functional dynamics

We can use the `nma()` function in bio3d to predict the large-scale functional motions of biomolecules

```
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
```
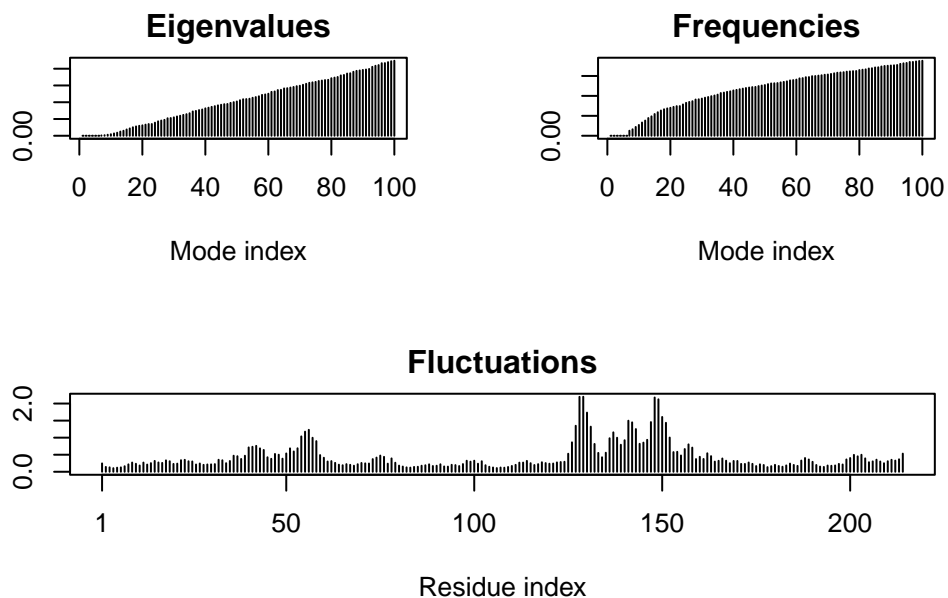
```
m <- nma(adk)
```

```
 Building Hessian...        Done in 0.073 seconds.
 Diagonalizing Hessian...   Done in 0.992 seconds.
```

```
plot(m)
```

12

## Eigenvalues

## Frequencies

## Fluctuations

Write out a trajectory of the predicted molecular motion:

```
mktrj(m, file="adk_m7.pdb")
```

```
library(bio3d)
aa <- get.seq("1ake_A")
```

Warning in get.seq("1ake_A"): Removing existing file: seqs.fasta

Fetching... Please wait. Done.

```
aa
```

```
            1        .         .         .         .         .        60
pdb|1AKE|A    MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMLRAAVKSGSELGKQAKDIMDAGKLVT
            1        .         .         .         .         .        60

            61       .         .         .         .         .        120
pdb|1AKE|A    DELVIALVKERIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFDVPDELIVDRI
            61       .         .         .         .         .        120
```

```
          121         .         .         .         .         .          180
pdb|1AKE|A   VGRRVHAPSGRVYHVKFNPPKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQMTAPLIG
          121         .         .         .         .         .          180

          181         .         .         .   214
pdb|1AKE|A   YYSKEAEAGNTKYAKVDGTKPVAEVRADLEKILG
          181         .         .         .   214
```

```
Call:
  read.fasta(file = outfile)

Class:
  fasta

Alignment dimensions:
  1 sequence rows; 214 position columns (214 non-gap, 0 gap)

+ attr: id, ali, call
```

```
b <- blast.pdb(aa)
```

```
 Searching ... please wait (updates every 5 seconds) RID = UG6FHMU9013
 .......................................................................
 Reporting 87 hits
```
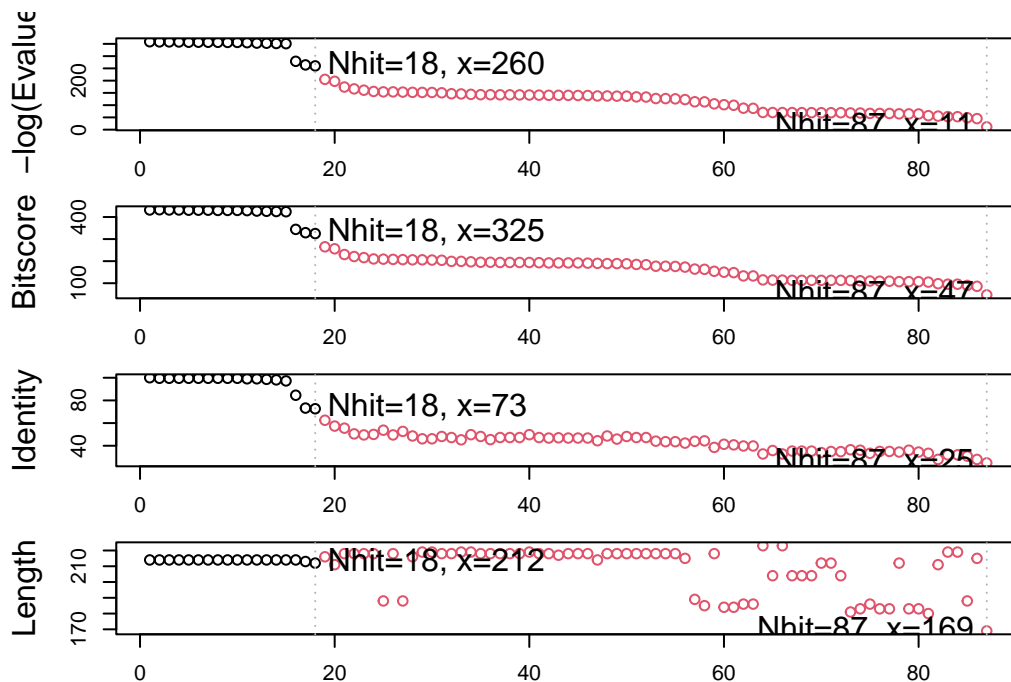
```
hits <- plot(b)
```

```
  * Possible cutoff values:    260 11
          Yielding Nhits:    18 87

  * Chosen cutoff value of:    260
          Yielding Nhits:    18
```

14

```
head(hits$pdb.id)
```

```
[1] "1AKE_A" "8BQF_A" "4X8M_A" "6S36_A" "8Q2B_A" "8RJ9_A"
```

```
hits <- NULL
hits$pdb.id <- c('1AKE_A','6S36_A','6RZE_A','3HPR_A','1E4V_A','5EJE_A','1E4Y_A','3X2S_A','6HA
files <- get.pdb(hits$pdb.id, path="pdbs", split=TRUE, gzip=TRUE)
```

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

```
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/6RZE.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/1E4V.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/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%
  |
  |=====                                                         |   8%
  |
  |==========                                                    |  15%
  |
  |===============                                               |  23%
  |
  |=====================                                         |  31%
  |
  |==========================                                    |  38%
```
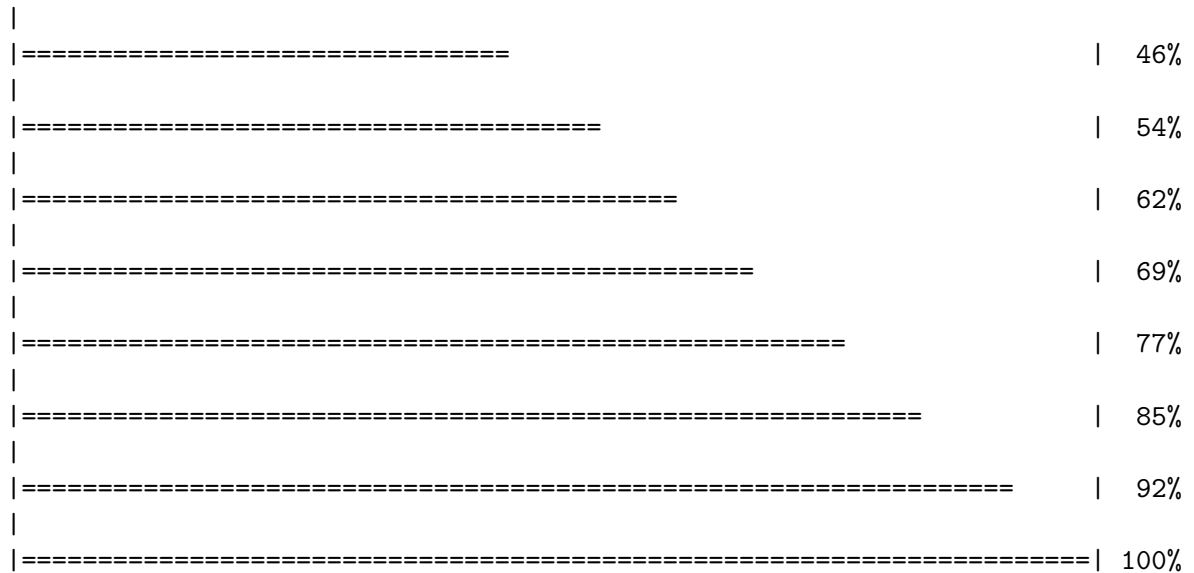
16

```
  |
  |==================================                                       |  46%
  |
  |=====================================                                    |  54%
  |
  |==========================================                               |  62%
  |
  |===============================================                          |  69%
  |
  |======================================================                   |  77%
  |
  |=============================================================            |  85%
  |
  |====================================================================     |  92%
  |
  |=========================================================================| 100%
```

```
pdbs <- pdbaln(files, fit = TRUE, exefile="msa")
```

```
Reading PDB files:
pdbs/split_chain/1AKE_A.pdb
pdbs/split_chain/6S36_A.pdb
pdbs/split_chain/6RZE_A.pdb
pdbs/split_chain/3HPR_A.pdb
pdbs/split_chain/1E4V_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/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
...
```

```
Extracting sequences

pdb/seq: 1   name: pdbs/split_chain/1AKE_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 2   name: pdbs/split_chain/6S36_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 3   name: pdbs/split_chain/6RZE_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 4   name: pdbs/split_chain/3HPR_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 5   name: pdbs/split_chain/1E4V_A.pdb
pdb/seq: 6   name: pdbs/split_chain/5EJE_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 7   name: pdbs/split_chain/1E4Y_A.pdb
pdb/seq: 8   name: pdbs/split_chain/3X2S_A.pdb
pdb/seq: 9   name: pdbs/split_chain/6HAP_A.pdb
pdb/seq: 10   name: pdbs/split_chain/6HAM_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 11   name: pdbs/split_chain/4K46_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 12   name: pdbs/split_chain/3GMT_A.pdb
pdb/seq: 13   name: pdbs/split_chain/4PZL_A.pdb
```
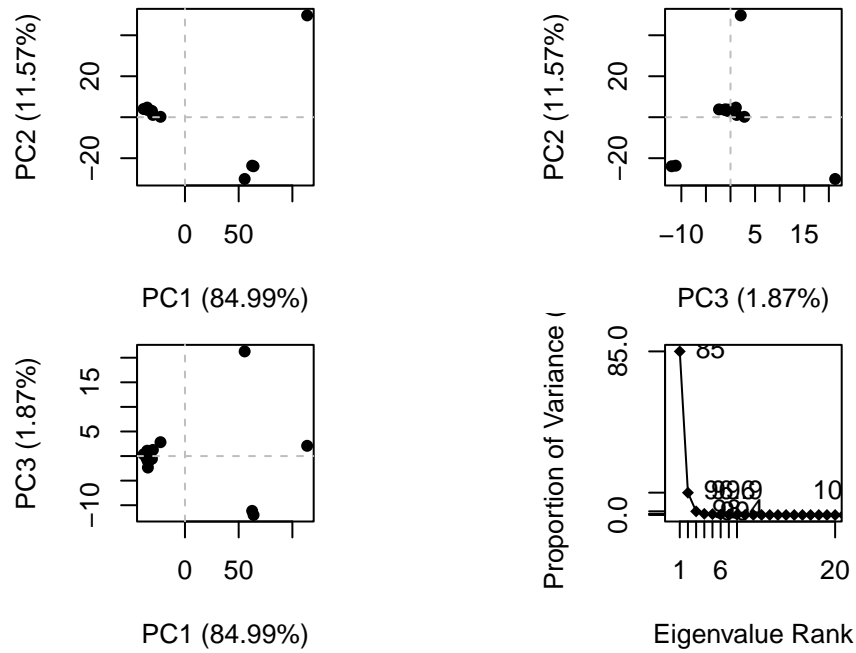
```r
ids <- basename.pdb(pdbs$id)
```

```r
anno <- pdb.annotate(ids)
unique(anno$source)
```

```
[1] "Escherichia coli"
[2] "Escherichia coli K-12"
[3] "Escherichia coli O139:H28 str. E24377A"
[4] "Escherichia coli str. K-12 substr. MDS42"
[5] "Photobacterium profundum"
[6] "Burkholderia pseudomallei 1710b"
[7] "Francisella tularensis subsp. tularensis SCHU S4"
```

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

```
rd <- rmsd(pdbs)
```

```
Warning in rmsd(pdbs): No indices provided, using the 204 non NA positions
```

```
hc.rd <- hclust(dist(rd))
grps.rd <- cutree(hc.rd, k=3)

plot(pc.xray, 1:2, col="grey50", bg=grps.rd, pch=21, cex=1)
```

19