# Class 9: Structural Bioinformatics

## Challana Tea

**What is in the PDB databse?**

```
pdbstats <- read.csv('pdb.csv', row.names = 1)
knitr::kable(pdbstats)
```

|                            | X.ray   | EM    | NMR    | Multiple.methods | Neutron | Other | Total   |
| -------------------------- | ------- | ----- | ------ | ---------------- | ------- | ----- | ------- |
| Protein (only)             | 152,809 | 9,421 | 12,117 | 191              | 72      | 32    | 174,642 |
| Protein/Oligosaccharide    | 9,008   | 1,654 | 32     | 7                | 1       | 0     | 10,702  |
| Protein/NA                 | 8,061   | 2,944 | 281    | 6                | 0       | 0     | 11,292  |
| Nucleic acid (only)        | 2,602   | 77    | 1,433  | 12               | 2       | 1     | 4,127   |
| Other                      | 163     | 9     | 31     | 0                | 0       | 0     | 203     |
| Oligosaccharide (only)     | 11      | 0     | 6      | 1                | 0       | 4     | 22      |

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

```
pdbstats$X.ray
```

```
[1] "152,809" "9,008"   "8,061"   "2,602"   "163"      "11"
```

```
as.numeric(pdbstats$X.ray)
```

```
Warning: NAs introduced by coercion
```

```
[1]  NA  NA  NA  NA 163  11
```

The commas are converting values containing them to NA values

```
n.xray <- sum(as.numeric(gsub(",","", pdbstats$X.ray)))
n.xray
```

[1] 172654

```
n.em <- sum(as.numeric(gsub(",","", pdbstats$EM)))
n.total <- sum(as.numeric(gsub(",","", pdbstats$Total)))

round(((n.xray + n.em)/n.total), 2)
```

[1] 0.93

```
rm_comma_sum <- function(filecolumn) {
  sum(as.numeric(gsub(",","", filecolumn)))
}

rm_comma_sum(pdbstats$X.ray)/rm_comma_sum(pdbstats$Total)
```

[1] 0.8590264

```
round(n.xray/n.total, 2)
```

[1] 0.86

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

```
round((rm_comma_sum(pdbstats$Total[1])/n.total), 2)
```

[1] 0.87

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

## Using the Molstar viewer

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

The structure has only been resolved to 2 A, which is above the resolution needed to see hydrogen.

> 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

H 308

> Q6: Generate and save a figure clearly showing the two distinct chains of HIV-protease along with the ligand. You might also consider showing the catalytic residues ASP 25 in each chain and the critical water (we recommend "Ball & Stick" for these side-chains). Add this figure to your Quarto document.

## Let's do some bioinformatics

Use the bio3D package for structural bioinformatics.

```
library(bio3d)
p <- read.pdb("1hsg")
```

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

```
p
```

```
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) ]
```
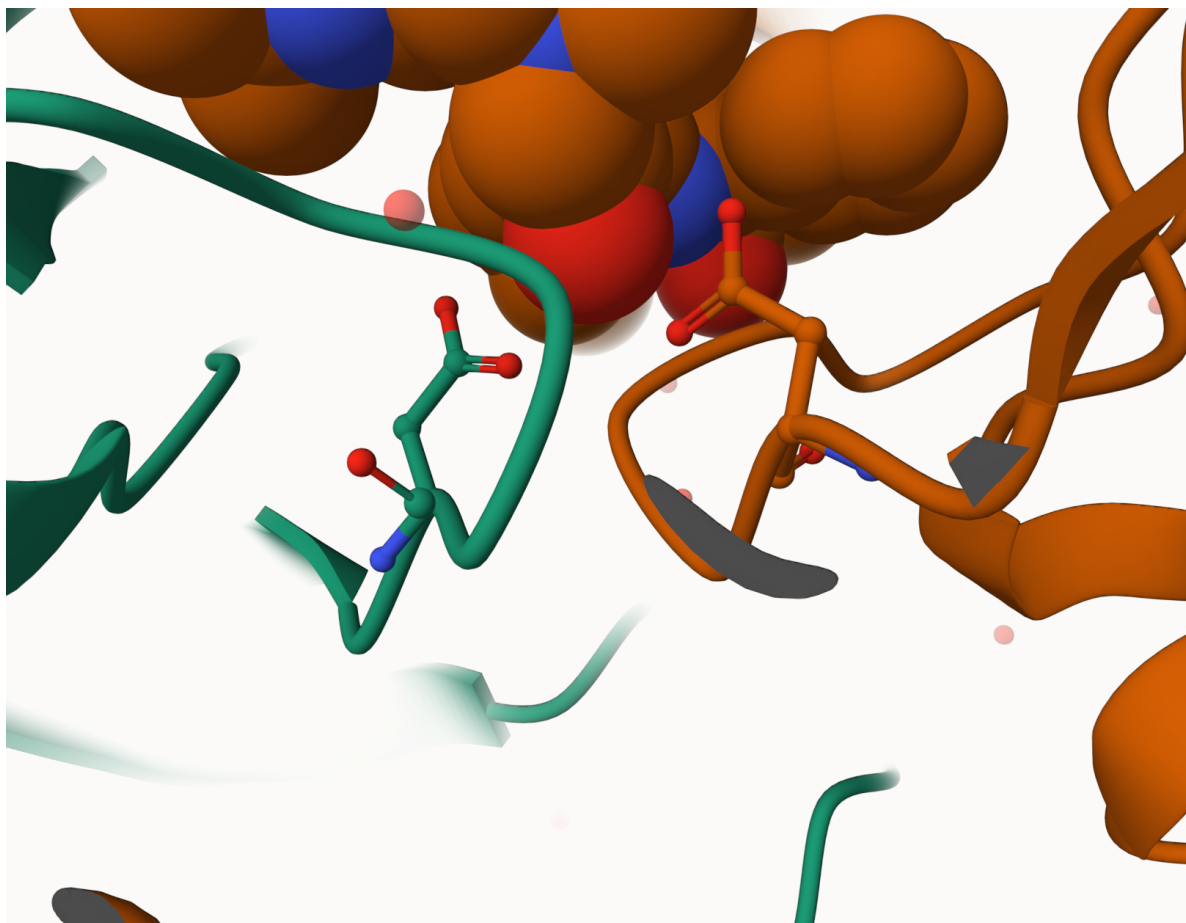
Figure 1: Fig. 1: A rendering of HIV-1 Pr active site with a bound ligand

```
  Protein sequence:
     PQITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWKPKMIGGIGGFIKVRQYD
     QILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNFPQITLWQRPLVTIKIGGQLKE
     ALLDTGADDTVLEEMSLPGRWKPKMIGGIGGFIKVRQYDQILIEICGHKAIGTVLVGPTP
     VNIIGRNLLTQIGCTLNF

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

```r
head(p$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>
```

```r
p$atom[1,'resid']
```

```
[1] "PRO"
```

```r
p$atom$resid[1]
```

```
[1] "PRO"
```

```r
aa321(p$atom$resid[1])
```

```
[1] "P"
```

## Normal Mode Analysis (NMA)

```
#Read an input structure
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
```
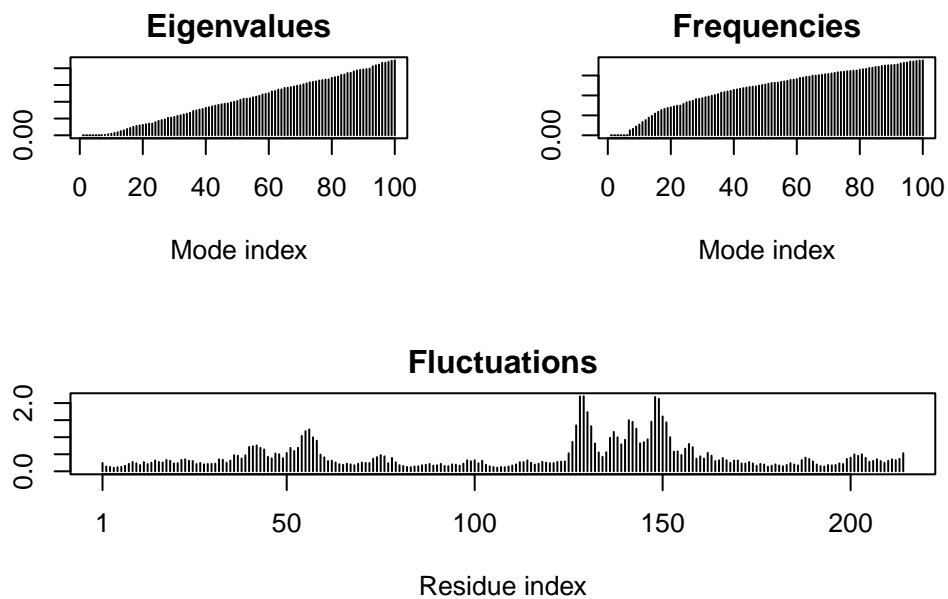
```
#NMA
m <- nma(adk)
```

```
Building Hessian...       Done in 0.028 seconds.
Diagonalizing Hessian...  Done in 0.371 seconds.
```

```
plot(m)
```

**Eigenvalues**

**Frequencies**

**Fluctuations**

Make a video of this motion for Molstar

```
#Make a trajectory file
mktrj(m, file="adk_m7.pdb")
```

# PCA - Comparative structure analysis of Adenylate Kinase

First, extract the sequence

```
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
```

```r
  #b <- blast.pdb(aa)
  #hits <- plot(b)
  # hits
```

```r
  hits <- NULL
  hits$pdb.id <- c('1AKE_A','6S36_A','6RZE_A','3HPR_A','1E4V_A','5EJE_A','1E4Y_A','3X2S_A','
```

```r
  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
```

8

```
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
```
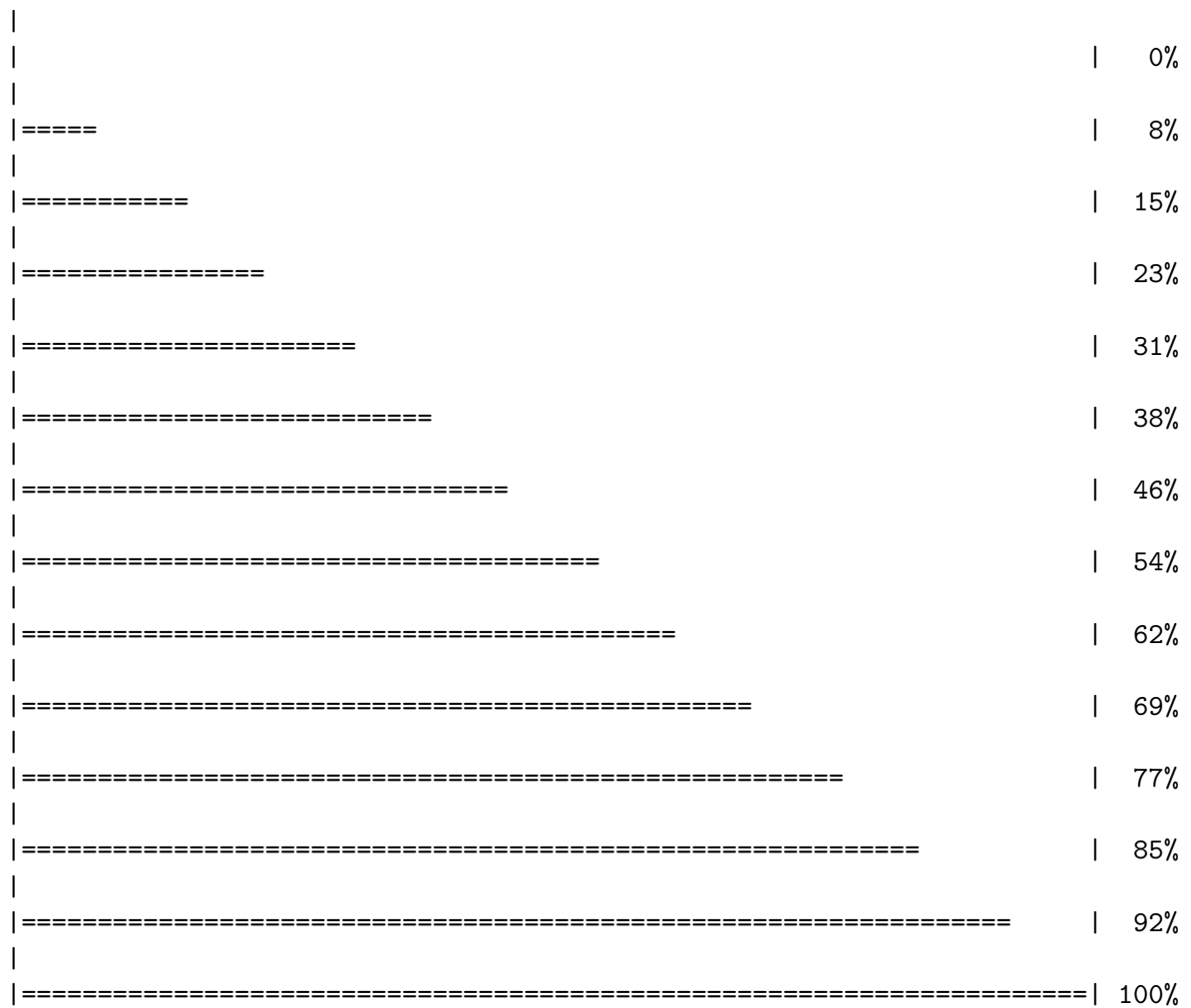
```
  |
  |                                                                     |   0%
  |=====                                                                |   8%
  |==========                                                           |  15%
  |===============                                                      |  23%
  |====================                                                 |  31%
  |=========================                                            |  38%
  |==============================                                       |  46%
  |===================================                                  |  54%
  |========================================                             |  62%
  |=============================================                        |  69%
  |==================================================                   |  77%
  |=======================================================              |  85%
  |============================================================         |  92%
  |=====================================================================| 100%
```

## Align and Superimpose

```r
#Align related pdbs
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
```
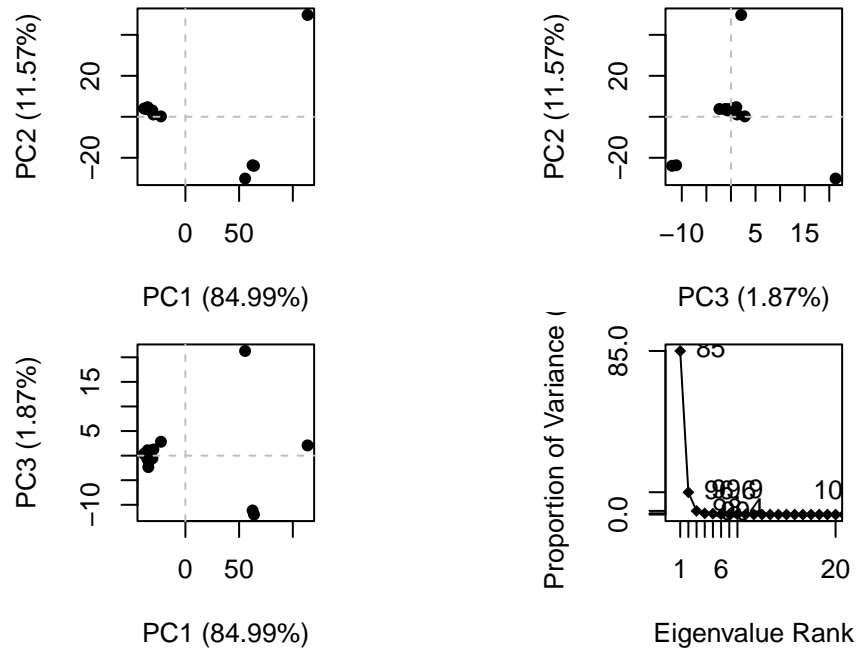
```
ids <- basename.pdb(pdbs$id)
graphics.off()
par(mar=c(0.5,0.5,0.5,0))
plot(pdbs, labels=ids, width = 5, height = 4)

#Run a PCA
pc.xray <- pca(pdbs)
plot(pc.xray)
```



Make a trajectory of the displacements captured by PCA

```
mktrj(pc.xray, pc = 1, file = "pc_1.pdb")
```