# Class 10: Structural Bioinformatics (Pt. 1)

Kris Price (PID: A17464127

## Table of contents

## Introduction to the RCSB Protein Data Bank (PDB)

```
pdb <- read.csv("pdb_stats.csv")
```

## PDB Statistics

Q1: What is the proportion of each method in the PDB?

```
sum(pdb$X.ray) / sum(pdb$Total)
```

```
[1] 0.8095077
```

```
sum(pdb$EM) / sum(pdb$Total)
```

```
[1] 0.1283843
```

```
sum(pdb$NMR) / sum(pdb$Total)
```

```
[1] 0.05902786
```

```
sum(pdb$Integrative) / sum(pdb$Total)
```

```
[1] 0.001534026
```

```
sum(pdb$Multiple.methods) / sum(pdb$Total)
```

```
[1] 0.001044101
```

```
sum(pdb$Neutron) / sum(pdb$Total)
```

```
[1] 0.0003533881
```

```
sum(pdb$Other) / sum(pdb$Total)
```

```
[1] 0.0001485836
```

- X-ray: 80.95%
- EM: 12.84%
- NMR: 5.9%
- Integrative: 0.15%
- Multiple methods: 0.1%
- Neutron: 0.04%
- Other: 0.01%

    Q2: What is the total number of entries in the PDB?

```
sum(pdb$Total)
```

```
[1] 249018
```

There are 249,018 entries in the PDB.

> 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?

4,940 HIV-1 protease structures are in the current PDB.

## Visualizing the HIV-1 protease structure

### Delving deeper

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

> 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.

> Discussion Topic: Can you think of a way in which indinavir, or even larger ligands
> and substrates, could enter the binding site?

## Introduction to Bio3D in R

### Reading PDB file data into 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?

198 amino acid residues.

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

HOH.

Q9: How many protein chains are in this structure?

There are 2 protein chains.

To find the attributes of any object, you can use the `attributes()` function:

```
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>
```

**Quick PDB visualization in R**

```
library(bio3dview)
library(NGLVieweR)

view.pdb(pdb) |>
  setSpin()
```

file:///C:/Users/kryan/AppData/Local/Temp/RtmpE1G4ZT/file5fc4579b3451/widget5fc426af4848.html

Let's custom color the chains and highlight some key residues as spacefill/vdw:

```r
# Select the important ASP 25 residue
sele <- atom.select(pdb, resno=25)

view.pdb(pdb, cols=c("navy","teal"),
         highlight = sele,
         highlight.style = "spacefill") |>
  setRock()
```

file:///C:/Users/kryan/AppData/Local/Temp/RtmpE1G4ZT/file5fc457f41683/widget5fc41f5c2570.html

**Predicting functional motions of a single structure**

```
adk <- read.pdb("6s36")
```
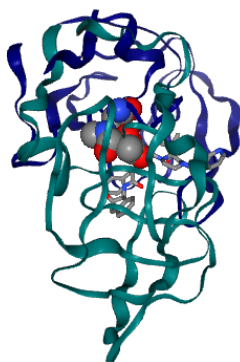
```
 Note: Accessing on-line PDB file
  PDB has ALT records, taking A only, rm.alt=TRUE
```

```
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
```

Normal mode analysis (NMA) is a structural bioinformatics method to predict protein flexibility and potential functional motions (a.k.a. conformational changes):

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

```
 Building Hessian...        Done in 0.03 seconds.
 Diagonalizing Hessian...   Done in 0.36 seconds.
```

```
plot(m)
```

**Eigenvalues**

**Frequencies**

**Fluctuations**

We can view a "movie" of these predicted motions by generating a molecular "trajectory" with the `mktrj()` function:

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

For a quicker display you can use the `view.nma()` function from the bio3dview package mentioned previously:

```
view.nma(m, pdb=adk)
```

```
file:///C:/Users/kryan/AppData/Local/Temp/RtmpE1G4ZT/file5fc479e6d5a/widget5fc42edf6a57.html
```

10

## Comparative structure analysis of Adenylate Kinase

### Setup

Q10. Which of the packages above is found only on BioConductor and not CRAN?

The `msa` package.

Q11. Which of the above packages is not found on BioConductor or CRAN?

The `bio3dview` package.

Q12. True or False? Functions from the pak package can be used to install packages from GitHub and BitBucket?

True.

### Search and retrieve ADK structures

```
library(bio3d)
aa <- get.seq("1ake_A")
```
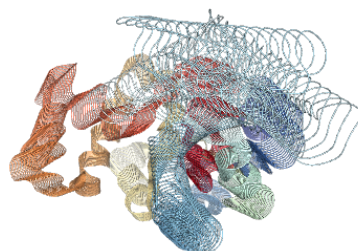
```
Warning in get.seq("1ake_A"): Removing existing file: seqs.fasta
```

```
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
```

Q13. How many amino acids are in this sequence, i.e. how long is this sequence?

214 amino acids.

**Align and superpose structures**

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

Using `get.pdb()` and `pdbslit()` functions to download and parse the above structures:

```
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 exists. Skipping download
```

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

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

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

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

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

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

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

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

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

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

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

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


  |
  |                                                                  |   0%
  |
  |=====                                                             |   8%
  |
  |==========                                                        |  15%
  |
  |===============                                                   |  23%
  |
  |=====================                                             |  31%
  |
  |=========================                                         |  38%
```

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

We can use the `pdbaln()` function to align and optionally fit (i.e. superpose) the PDB structures.
Then, we could use `pdb.annotate()` to annotate each structure to its source species:

```
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)

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
anno
```

```
      structureId chainId macromoleculeType chainLength experimentalTechnique
```

| | | | | | |
|---|---|---|---|---|---|
| 1AKE_A | 1AKE | A | Protein | 214 | X-ray |
| 6S36_A | 6S36 | A | Protein | 214 | X-ray |
| 6RZE_A | 6RZE | A | Protein | 214 | X-ray |
| 3HPR_A | 3HPR | A | Protein | 214 | X-ray |
| 1E4V_A | 1E4V | A | Protein | 214 | X-ray |
| 5EJE_A | 5EJE | A | Protein | 214 | X-ray |
| 1E4Y_A | 1E4Y | A | Protein | 214 | X-ray |
| 3X2S_A | 3X2S | A | Protein | 214 | X-ray |
| 6HAP_A | 6HAP | A | Protein | 214 | X-ray |
| 6HAM_A | 6HAM | A | Protein | 214 | X-ray |
| 4K46_A | 4K46 | A | Protein | 214 | X-ray |
| 3GMT_A | 3GMT | A | Protein | 230 | X-ray |
| 4PZL_A | 4PZL | A | Protein | 242 | X-ray |

| | resolution | scopDomain | pfam |
|---|---|---|---|
| 1AKE_A | 2.00 | Adenylate kinase | Adenylate kinase, active site lid (ADK_lid) |
| 6S36_A | 1.60 | <NA> | Adenylate kinase (ADK) |
| 6RZE_A | 1.69 | <NA> | Adenylate kinase, active site lid (ADK_lid) |
| 3HPR_A | 2.00 | <NA> | Adenylate kinase (ADK) |
| 1E4V_A | 1.85 | Adenylate kinase | Adenylate kinase (ADK) |
| 5EJE_A | 1.90 | <NA> | Adenylate kinase, active site lid (ADK_lid) |
| 1E4Y_A | 1.85 | Adenylate kinase | Adenylate kinase (ADK) |
| 3X2S_A | 2.80 | <NA> | <NA> |
| 6HAP_A | 2.70 | <NA> | Adenylate kinase, active site lid (ADK_lid) |
| 6HAM_A | 2.55 | <NA> | Adenylate kinase (ADK) |
| 4K46_A | 2.01 | <NA> | Adenylate kinase, active site lid (ADK_lid) |
| 3GMT_A | 2.10 | <NA> | <NA> |
| 4PZL_A | 2.10 | <NA> | Adenylate kinase, active site lid (ADK_lid) |

| | ligandId |
|---|---|
| 1AKE_A | AP5 |
| 6S36_A | MG (2),NA,CL (3) |

```
6RZE_A      NA (3),CL (2)
3HPR_A              AP5
1E4V_A              AP5
5EJE_A           AP5,CO
1E4Y_A              AP5
3X2S_A    AP5,MG,JPY (2)
6HAP_A              AP5
6HAM_A              AP5
4K46_A      ADP,PO4,AMP
3GMT_A           SO4 (2)
4PZL_A       CA,FMT,GOL
                                                                ligandName
1AKE_A                                              BIS(ADENOSINE)-5'-
PENTAPHOSPHATE
6S36_A                          MAGNESIUM ION (2),SODIUM ION,CHLORIDE ION (3)
6RZE_A                                    SODIUM ION (3),CHLORIDE ION (2)
3HPR_A                                              BIS(ADENOSINE)-5'-
PENTAPHOSPHATE
1E4V_A                                              BIS(ADENOSINE)-5'-
PENTAPHOSPHATE
5EJE_A                          BIS(ADENOSINE)-5'-PENTAPHOSPHATE,COBALT (II) ION
1E4Y_A                                              BIS(ADENOSINE)-5'-
PENTAPHOSPHATE
3X2S_A BIS(ADENOSINE)-5'-PENTAPHOSPHATE,MAGNESIUM ION,N-(pyren-1-ylmethyl)acetamide (2)
6HAP_A                                              BIS(ADENOSINE)-5'-
PENTAPHOSPHATE
6HAM_A                                              BIS(ADENOSINE)-5'-
PENTAPHOSPHATE
4K46_A                ADENOSINE-5'-DIPHOSPHATE,PHOSPHATE ION,ADENOSINE MONOPHOSPHATE
3GMT_A                                                     SULFATE ION (2)
4PZL_A                                    CALCIUM ION,FORMIC ACID,GLYCEROL
                                                source
1AKE_A                           Escherichia coli
6S36_A                           Escherichia coli
6RZE_A                           Escherichia coli
3HPR_A                      Escherichia coli K-12
1E4V_A                           Escherichia coli
5EJE_A         Escherichia coli O139:H28 str. E24377A
1E4Y_A                           Escherichia coli
3X2S_A      Escherichia coli str. K-12 substr. MDS42
6HAP_A         Escherichia coli O139:H28 str. E24377A
6HAM_A                      Escherichia coli K-12
4K46_A                 Photobacterium profundum
```

```
3GMT_A                    Burkholderia pseudomallei 1710b
4PZL_A Francisella tularensis subsp. tularensis SCHU S4


1AKE_A STRUCTURE OF THE COMPLEX BETWEEN ADENYLATE KINASE FROM ESCHERICHIA COLI AND THE INHIBI
6S36_A
6RZE_A
3HPR_A
1E4V_A
loop
5EJE_A                                                                          Cryst
1E4Y_A
loop
3X2S_A
conjugated adenylate kinase
6HAP_A
6HAM_A
4K46_A
3GMT_A
4PZL_A                                                                      The cryst
```
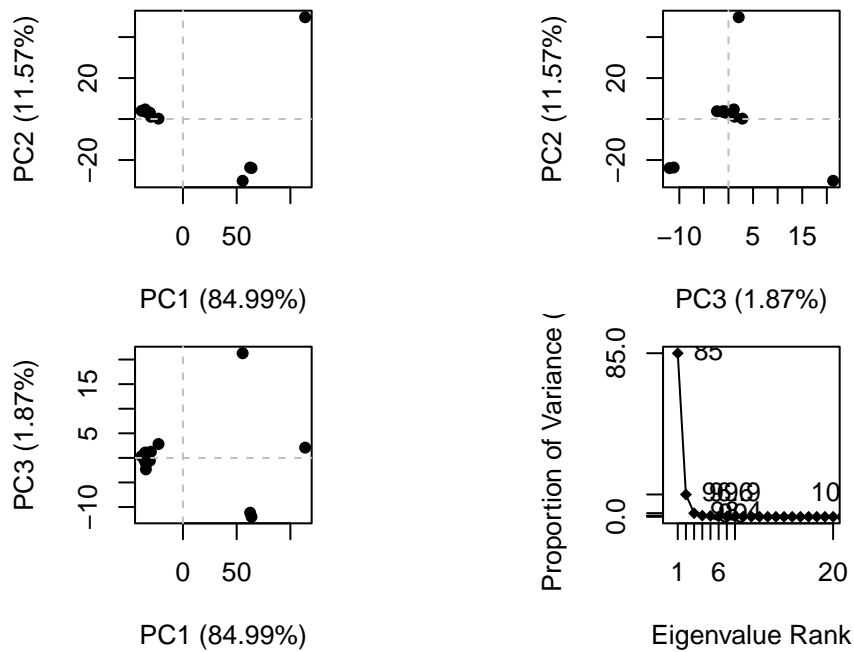
|        | citation | rObserved | rFree |
|--------|----------|-----------|-------|
| 1AKE_A | Muller, C.W., et al. J Mol Biology (1992) | 0.19600 | NA |
| 6S36_A | Rogne, P., et al. Biochemistry (2019) | 0.16320 | 0.23560 |
| 6RZE_A | Rogne, P., et al. Biochemistry (2019) | 0.18650 | 0.23500 |
| 3HPR_A | Schrank, T.P., et al. Proc Natl Acad Sci U S A (2009) | 0.21000 | 0.24320 |
| 1E4V_A | Muller, C.W., et al. Proteins (1993) | 0.19600 | NA |
| 5EJE_A | Kovermann, M., et al. Proc Natl Acad Sci U S A (2017) | 0.18890 | 0.23580 |
| 1E4Y_A | Muller, C.W., et al. Proteins (1993) | 0.17800 | NA |
| 3X2S_A | Fujii, A., et al. Bioconjug Chem (2015) | 0.20700 | 0.25600 |
| 6HAP_A | Kantaev, R., et al. J Phys Chem B (2018) | 0.22630 | 0.27760 |
| 6HAM_A | Kantaev, R., et al. J Phys Chem B (2018) | 0.20511 | 0.24325 |
| 4K46_A | Cho, Y.-J., et al. To be published | 0.17000 | 0.22290 |
| 3GMT_A | Buchko, G.W., et al. Biochem Biophys Res Commun (2010) | 0.23800 | 0.29500 |
| 4PZL_A | Tan, K., et al. To be published | 0.19360 | 0.23680 |

|        | rWork | spaceGroup |
|--------|-------|------------|
| 1AKE_A | 0.19600 | P 21 2 21 |
| 6S36_A | 0.15940 | C 1 2 1 |
| 6RZE_A | 0.18190 | C 1 2 1 |
| 3HPR_A | 0.20620 | P 21 21 2 |
| 1E4V_A | 0.19600 | P 21 2 21 |
| 5EJE_A | 0.18630 | P 21 2 21 |
| 1E4Y_A | 0.17800 | P 1 21 1 |
| 3X2S_A | 0.20700 | P 21 21 21 |
| 6HAP_A | 0.22370 | I 2 2 2 |

```
6HAM_A 0.20311      P 43
4K46_A 0.16730 P 21 21 21
3GMT_A 0.23500   P 1 21 1
4PZL_A 0.19130      P 32
```

**Principle component analysis**

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



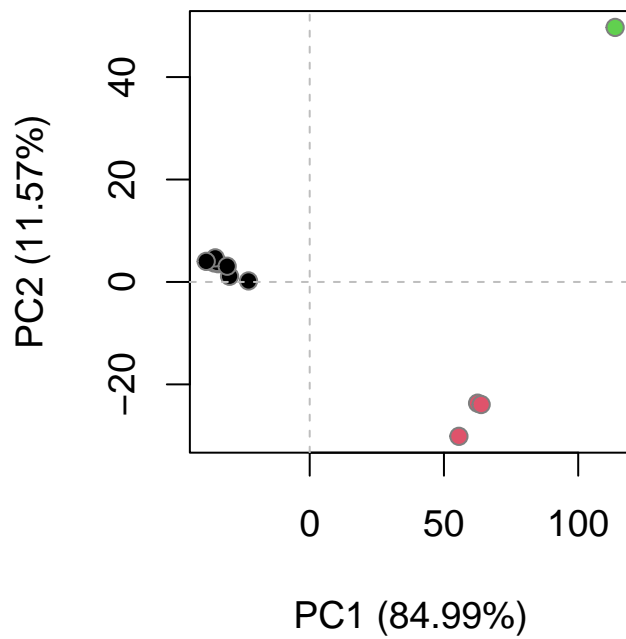The `rmsd()` function will calculate all pairwise RMSD values of thet structural ensemble, which can facilitate clustering analysis based on the pairwise structural deviation:

```
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)
```

**Normal Mode Analysis**

```
modes <- nma(pdbs)
```

```
Details of Scheduled Calculation:
  ... 13 input structures
  ... storing 606 eigenvectors for each structure
  ... dimension of x$U.subspace: ( 612x606x13 )
  ... coordinate superposition prior to NM calculation
  ... aligned eigenvectors (gap containing positions removed)
  ... estimated memory usage of final 'eNMA' object: 36.9 Mb


|
|                                                                |   0%
|
|=====                                                           |   8%
|
|==========                                                      |  15%
|
```

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

```r
plot(modes, pdbs, col=grps.rd)
```

Extracting SSE from pdbs$sse attribute

Q14. What do you note about this plot? Are the black and colored lines similar or different? Where do you think they differ most and why?

The black and colored lines both have similar peaks, but the colored lines have much greater fluctuations.