# class09

Joshua Lau

## What is in the PDB?

PDB is the main database of biomolecular structures, and is available at rcsb.org

First we load the dataset into R

```
PDB <- read.csv("./PDB.csv", row.names = 1)
#PDB <- apply(PDB, 2, as.numeric)
```

Looks like there are some problems with the numbers saved in the PDB file. Firstly, they are saved as characters. Secondly, they have commas in them. Therefore, we define a function to eliminate commas from a vector of strings.

```
replaceCommas<-function(x){
  return(as.numeric(gsub("\\,", "", x)))
}
```

Let's try to use this on the first column of the PDB dataset

```
replaceCommas(PDB$X.ray)
```

```
[1] 152809   9008   8061   2602    163     11
```

Looks like it works! Our next step is to apply this to all columns of the PDB dataset

```
PDB_num <- as.data.frame(apply(PDB, 2, replaceCommas))
rownames(PDB_num) <- rownames(PDB)
```

Now we can answer our question

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

```
sum(PDB_num$X.ray) / sum(PDB_num$Total) * 100
```

[1] 85.90264

```
sum(PDB_num$EM) / sum(PDB_num$Total) * 100
```

[1] 7.017832

85.9% of structures in the PDB are solved by X-ray and 7.0% of structures by EM

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

```
PDB_num$Total[1] / sum(PDB_num$Total) * 100
```

[1] 86.89175

86.9% of structures in the PDB are protein

**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?**

It is not straightforward to find all HIV-1 protease structures using plain text searching on databases.

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

We only see one atom per water molecule in this structure because the crystal structure does not have the sufficient resolution to view/resolve Hydrogen atoms, as they are too small in size.

**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**

The water molecule has the residue number 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.**



# Intro to bio3d in R

```
library(bio3d)
```

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

Loading the 1hsg structure

```
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

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

HOH

## Q9: How many protein chains are in this structure?

2

```
  attributes(pdb)
```

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


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

## Predicting functional motions of a single structure

Let's read a new PDB structure of Adenylate Kinase and perform Normal mode analysis.

```
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
```
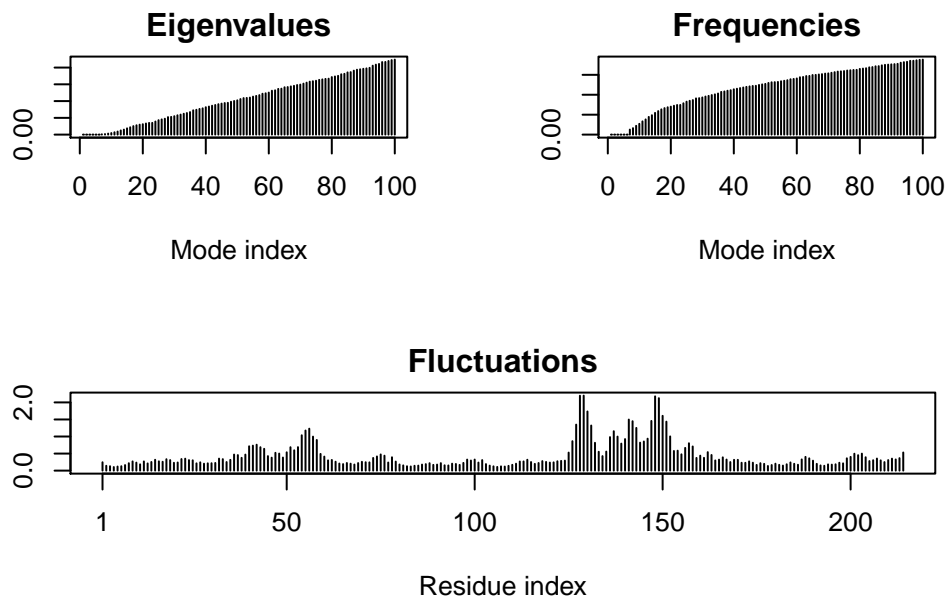
```
# Perform flexiblity prediction
m <- nma(adk)
```

```
Building Hessian...       Done in 0.02 seconds.
Diagonalizing Hessian...   Done in 0.41 seconds.
```

```
plot(m)
```

**Eigenvalues**

**Frequencies**

**Fluctuations**

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

# Section 4. Comparative structure analysis

```
library(bio3d)
library(BiocManager)
```

```
Bioconductor version '3.12' is out-of-date; the current release version '3.16'
  is available with R version '4.2'; see https://bioconductor.org/install
```

**Q10.**

the msa package

**Q11.**

bio3d-view

**Q12.**

TRUE

## Search and retrieve ADK structures

```
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.

**214**

```r
#b <- blast.pdb(aa)
#we don't want to run this every time

#saveRDS(b, file = "blast_results.RDS")
```

```r
b <- readRDS("blast_results.RDS")
```

```r
hits <- plot(b)
```

```
* Possible cutoff values:    197 -3
         Yielding Nhits:    16 96

* Chosen cutoff value of:    197
         Yielding Nhits:    16
```

```r
hits$pdb.id
```

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

```r
# 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/
1AKE.pdb exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE): pdbs/
4X8M.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/
4X8H.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/
4NP6.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
```
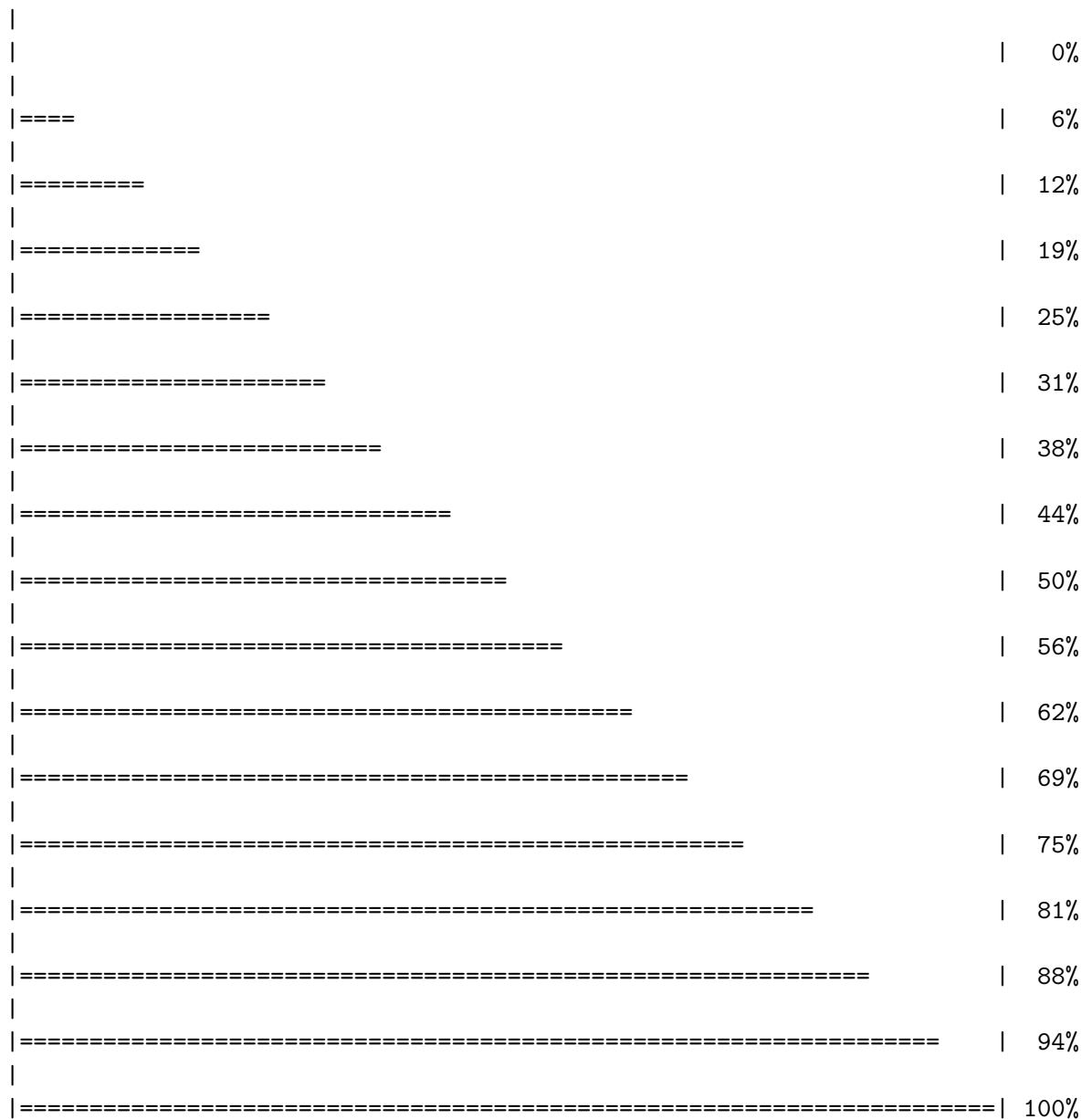
```
|
|                                                                      |   0%
|
|====                                                                  |   6%
|
|========                                                              |  12%
|
|============                                                          |  19%
|
|=================                                                     |  25%
|
|=====================                                                 |  31%
|
|=========================                                             |  38%
|
|=============================                                         |  44%
|
|=================================                                     |  50%
|
|=====================================                                 |  56%
|
|=========================================                             |  62%
|
|=============================================                         |  69%
|
|==================================================                    |  75%
|
|======================================================                |  81%
|
|==========================================================            |  88%
|
|==============================================================        |  94%
|
|======================================================================| 100%
```

Next we are going to align and superimpose all these structures

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

```
Reading PDB files:
pdbs/split_chain/1AKE_A.pdb
```

```
pdbs/split_chain/4X8M_A.pdb
pdbs/split_chain/6S36_A.pdb
pdbs/split_chain/6RZE_A.pdb
pdbs/split_chain/4X8H_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/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
....

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/4X8M_A.pdb
pdb/seq: 3   name: pdbs/split_chain/6S36_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/4X8H_A.pdb
pdb/seq: 6   name: pdbs/split_chain/3HPR_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 7   name: pdbs/split_chain/1E4V_A.pdb
pdb/seq: 8   name: pdbs/split_chain/5EJE_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 9   name: pdbs/split_chain/1E4Y_A.pdb
pdb/seq: 10   name: pdbs/split_chain/3X2S_A.pdb
pdb/seq: 11   name: pdbs/split_chain/6HAP_A.pdb
pdb/seq: 12   name: pdbs/split_chain/6HAM_A.pdb
```

```
    PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 13   name: pdbs/split_chain/4K46_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 14   name: pdbs/split_chain/4NP6_A.pdb
pdb/seq: 15   name: pdbs/split_chain/3GMT_A.pdb
pdb/seq: 16   name: pdbs/split_chain/4PZL_A.pdb
```

```r
  # Vector containing PDB codes for figure axis
  ids <- basename.pdb(pdbs$id)

  # Draw schematic alignment
  #plot(pdbs, labels=ids)


  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] "Vibrio cholerae O1 biovar El Tor str. N16961"
[7] "Burkholderia pseudomallei 1710b"
[8] "Francisella tularensis subsp. tularensis SCHU S4"
```

```r
  anno
```

| | structureId | chainId | macromoleculeType | chainLength | experimentalTechnique |
|---|---|---|---|---|---|
| 1AKE_A | 1AKE | A | Protein | 214 | X-ray |
| 4X8M_A | 4X8M | A | Protein | 214 | X-ray |
| 6S36_A | 6S36 | A | Protein | 214 | X-ray |
| 6RZE_A | 6RZE | A | Protein | 214 | X-ray |
| 4X8H_A | 4X8H | 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
4NP6_A          4NP6         A            Protein            217                X-ray
3GMT_A          3GMT         A            Protein            230                X-ray
4PZL_A          4PZL         A            Protein            242                X-ray
          resolution       scopDomain                        pfam           ligandId
1AKE_A       2.000 Adenylate kinase Adenylate kinase (ADK)               AP5
4X8M_A       2.600            <NA> Adenylate kinase (ADK)              <NA>
6S36_A       1.600            <NA> Adenylate kinase (ADK) CL (3),NA,MG (2)
6RZE_A       1.690            <NA> Adenylate kinase (ADK)    NA (3),CL (2)
4X8H_A       2.500            <NA> Adenylate kinase (ADK)              <NA>
3HPR_A       2.000            <NA> Adenylate kinase (ADK)               AP5
1E4V_A       1.850 Adenylate kinase Adenylate kinase (ADK)               AP5
5EJE_A       1.900            <NA> Adenylate kinase (ADK)           AP5,CO
1E4Y_A       1.850 Adenylate kinase Adenylate kinase (ADK)               AP5
3X2S_A       2.800            <NA> Adenylate kinase (ADK)   JPY (2),AP5,MG
6HAP_A       2.700            <NA> Adenylate kinase (ADK)               AP5
6HAM_A       2.550            <NA> Adenylate kinase (ADK)               AP5
4K46_A       2.010            <NA> Adenylate kinase (ADK)      ADP,AMP,PO4
4NP6_A       2.004            <NA> Adenylate kinase (ADK)              <NA>
3GMT_A       2.100            <NA> Adenylate kinase (ADK)           SO4 (2)
4PZL_A       2.100            <NA> Adenylate kinase (ADK)       GOL,CA,FMT
                                                                    ligandName
1AKE_A                                    BIS(ADENOSINE)-5'-PENTAPHOSPHATE
4X8M_A                                                              <NA>
6S36_A                     CHLORIDE ION (3),SODIUM ION,MAGNESIUM ION (2)
6RZE_A                              SODIUM ION (3),CHLORIDE ION (2)
4X8H_A                                                              <NA>
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 N-(pyren-1-ylmethyl)acetamide (2),BIS(ADENOSINE)-5'-PENTAPHOSPHATE,MAGNESIUM ION
6HAP_A                                    BIS(ADENOSINE)-5'-PENTAPHOSPHATE
6HAM_A                                    BIS(ADENOSINE)-5'-PENTAPHOSPHATE
4K46_A          ADENOSINE-5'-DIPHOSPHATE,ADENOSINE MONOPHOSPHATE,PHOSPHATE ION
4NP6_A                                                              <NA>
3GMT_A                                                   SULFATE ION (2)
4PZL_A                                GLYCEROL,CALCIUM ION,FORMIC ACID
                                        source
1AKE_A                       Escherichia coli
4X8M_A                       Escherichia coli
6S36_A                       Escherichia coli
6RZE_A                       Escherichia coli
```

14

```
4X8H_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
4NP6_A      Vibrio cholerae O1 biovar El Tor str. N16961
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 INHIB
4X8M_A
6S36_A
6RZE_A
4X8H_A
3HPR_A
1E4V_A
5EJE_A                                                                          Crys
1E4Y_A
3X2S_A
6HAP_A
6HAM_A
4K46_A
4NP6_A
3GMT_A
4PZL_A                                                                          The cryst
```
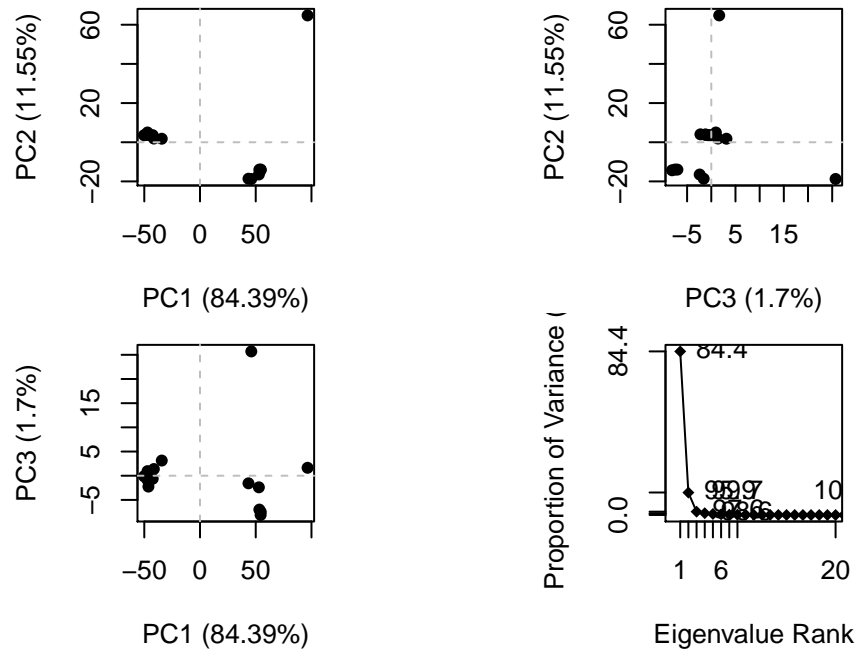
| | citation | rObserved | rFree |
|---|---|---|---|
| 1AKE_A | Muller, C.W., et al. J Mol Biol (1992) | 0.19600 | NA |
| 4X8M_A | Kovermann, M., et al. Nat Commun (2015) | 0.24910 | 0.30890 |
| 6S36_A | Rogne, P., et al. Biochemistry (2019) | 0.16320 | 0.23560 |
| 6RZE_A | Rogne, P., et al. Biochemistry (2019) | 0.18650 | 0.23500 |
| 4X8H_A | Kovermann, M., et al. Nat Commun (2015) | 0.19610 | 0.28950 |
| 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 |

```
4NP6_A                            Kim, Y., et al. To be published    0.18800 0.22200
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
4X8M_A 0.24630     C 1 2 1
6S36_A 0.15940     C 1 2 1
6RZE_A 0.18190     C 1 2 1
4X8H_A 0.19140     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
4NP6_A 0.18600        P 43
3GMT_A 0.23500    P 1 21 1
4PZL_A 0.19130        P 32
```

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

```
# Calculate RMSD
rd <- rmsd(pdbs)
```

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

```
# Structure-based clustering
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)
```