

Class 11

In the first section, we will interact with the main US based PDB website

First we will export the csv data from the website

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

	X.ray <int>	NMR <int>	EM <int>	Multiple.methods <int>	Neutr... <int>	Ot... <int>	T <
Protein (only)	142303	11804	5999	177	70	32	160
Protein/Oligosaccharide	8414	31	979	5	0	0	9
Protein/NA	7491	274	1986	3	0	0	9
Nucleic acid (only)	2368	1372	60	8	2	1	3
Other	149	31	3	0	0	0	
Oligosaccharide (only)	11	6	0	1	0	4	

6 rows

```
View(db)
```

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

```
method.sums <- colSums(db)
round(method.sums/method.sums["Total"] * 100, 2)
```

```
##          X.ray          NMR          EM Multiple.methods
##          87.55          7.36          4.92          0.11
##       Neutron          Other          Total
##          0.04          0.02         100.00
```

87.55% of structures in the PDB are solved of X-ray and 4.92% of structures in the PDB are solved by Electron Microscopy.

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

```
round(db$Total/method.sums["Total"] * 100,2 )
```

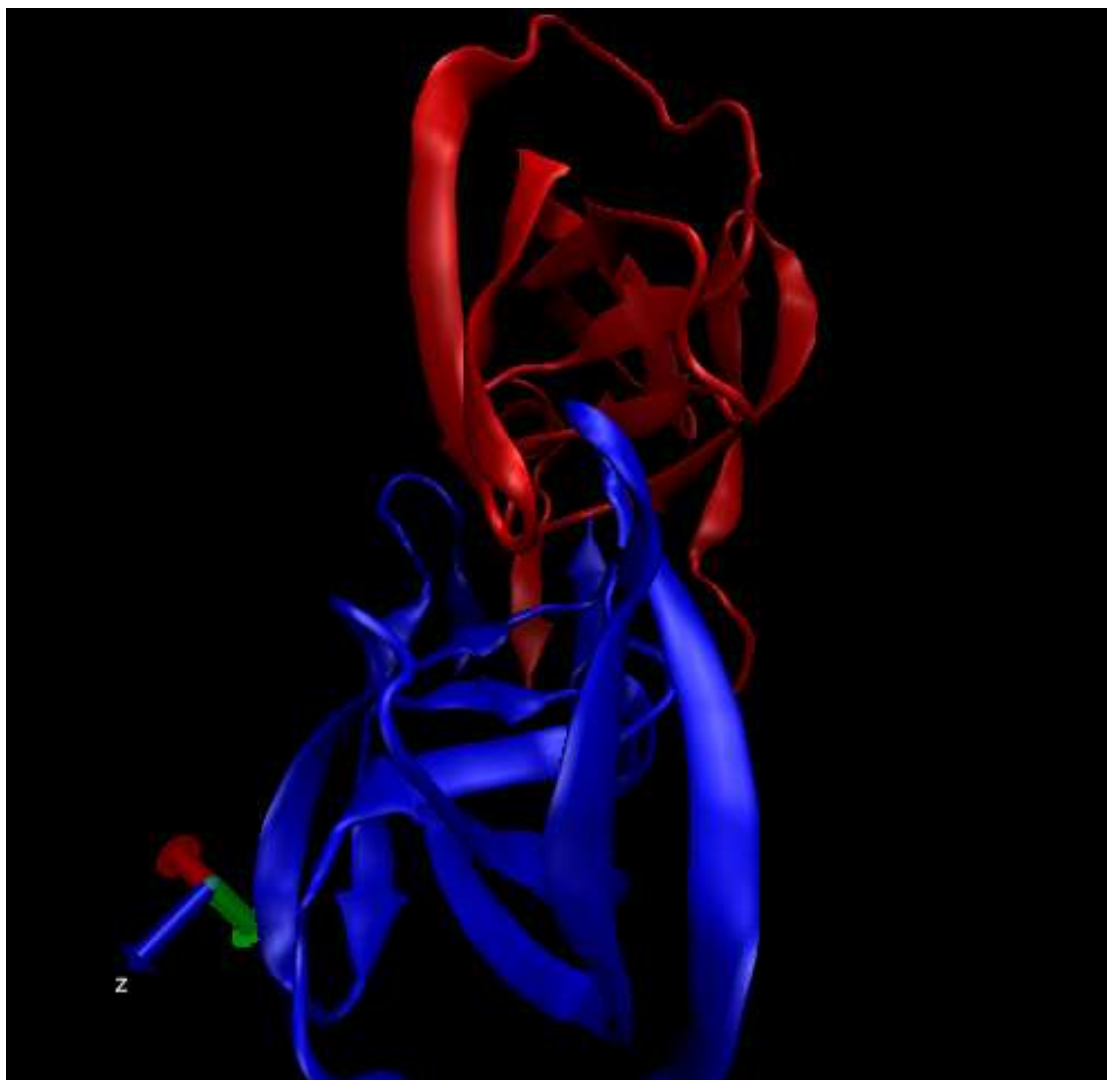
```
## [1] 87.36  5.14  5.31  2.08  0.10  0.01
```

87.36% 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?

1828 structures

VMD structure Visualization Image



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

It was selected to show not water and not protein.

Q5: There is a conserved water molecule in the binding site. Can you identify this water molecule? What residue number does this water molecule have (see note below)?

Using Bio3D

```
library(bio3d)

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

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

Q6 was in the optional section Q7: How many amino acid residues are there in this pdb object?

There are 198 amino acid residues.

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

MK1 is one of the two non-protein residues.

```
attributes(pdb)
```

```
## $names
## [1] "atom" "xyz" "seqres" "helix" "sheet" "calpha" "remark" "call"
##
## $class
## [1] "pdb" "sse"
```

Q9: How many protein chains are in this structure?

There are two protein chains.

```
#viewing full sequence
pdbseq(pdb)
```

```
##  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"
```

```
#viewing amino acid seq
aa123(pdbseq(pdb))
```

```
## [1] "PRO" "GLN" "ILE" "THR" "LEU" "TRP" "GLN" "ARG" "PRO" "LEU" "VAL" "THR"
## [13] "ILE" "LYS" "ILE" "GLY" "GLY" "GLN" "LEU" "LYS" "GLU" "ALA" "LEU" "LEU"
## [25] "ASP" "THR" "GLY" "ALA" "ASP" "ASP" "THR" "VAL" "LEU" "GLU" "GLU" "MET"
## [37] "SER" "LEU" "PRO" "GLY" "ARG" "TRP" "LYS" "PRO" "LYS" "MET" "ILE" "GLY"
## [49] "GLY" "ILE" "GLY" "GLY" "PHE" "ILE" "LYS" "VAL" "ARG" "GLN" "TYR" "ASP"
## [61] "GLN" "ILE" "LEU" "ILE" "GLU" "ILE" "CYS" "GLY" "HIS" "LYS" "ALA" "ILE"
## [73] "GLY" "THR" "VAL" "LEU" "VAL" "GLY" "PRO" "THR" "PRO" "VAL" "ASN" "ILE"
## [85] "ILE" "GLY" "ARG" "ASN" "LEU" "LEU" "THR" "GLN" "ILE" "GLY" "CYS" "THR"
## [97] "LEU" "ASN" "PHE" "PRO" "GLN" "ILE" "THR" "LEU" "TRP" "GLN" "ARG" "PRO"
## [109] "LEU" "VAL" "THR" "ILE" "LYS" "ILE" "GLY" "GLY" "GLN" "LEU" "LYS" "GLU"
## [121] "ALA" "LEU" "LEU" "ASP" "THR" "GLY" "ALA" "ASP" "ASP" "THR" "VAL" "LEU"
## [133] "GLU" "GLU" "MET" "SER" "LEU" "PRO" "GLY" "ARG" "TRP" "LYS" "PRO" "LYS"
## [145] "MET" "ILE" "GLY" "GLY" "ILE" "GLY" "GLY" "PHE" "ILE" "LYS" "VAL" "ARG"
## [157] "GLN" "TYR" "ASP" "GLN" "ILE" "LEU" "ILE" "GLU" "ILE" "CYS" "GLY" "HIS"
## [169] "LYS" "ALA" "ILE" "GLY" "THR" "VAL" "LEU" "VAL" "GLY" "PRO" "THR" "PRO"
## [181] "VAL" "ASN" "ILE" "ILE" "GLY" "ARG" "ASN" "LEU" "LEU" "THR" "GLN" "ILE"
## [193] "GLY" "CYS" "THR" "LEU" "ASN" "PHE"
```

The ATOM records

```
head(pdb$atom)
```

type <chr>	eleno <int>	elety <chr>	alt <chr>	resid <chr>	chain <chr>	resno <int>	insert <chr>	x <dbl>
1 ATOM	1	N	NA	PRO	A	1	NA	29.361
2 ATOM	2	CA	NA	PRO	A	1	NA	30.307
3 ATOM	3	C	NA	PRO	A	1	NA	29.760
4 ATOM	4	O	NA	PRO	A	1	NA	28.600
5 ATOM	5	CB	NA	PRO	A	1	NA	30.508
6 ATOM	6	CG	NA	PRO	A	1	NA	29.296

6 rows | 1-10 of 17 columns

Comparative structure analysis of Adenylate Kinase

The goal of this section is to perform PCA on the complete collection of Adenylate kinase structures in the PDB

```
# Install packages in the R console not your Rmd

#install.packages("bio3d")
#install.packages("ggplot2")
#install.packages("ggrepel")
#install.packages("devtools")
#install.packages("BiocManager")

#BiocManager::install("msa")
#devtools::install_bitbucket("Grantlab/bio3d-view")
```

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

msa is found only on BioConductor and not CRAN

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

bio3d-view

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

TRUE

Search and retrieve ADK structures

We will perform a blast search of the PDB database to identify related structures to our ADK sequence

```
library(bio3d)
# fetch the query sequence for chain A of the PDB ID 1AKE
aa <- get.seq("1ake_A")
```

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

```
## Fetching... Please wait. Done.
```

```
#viewing query sequence
aa
```

```

##          1          .          .          .          .          .          60
## pdb|1AKE|A  MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMLRAAVKSGSELGKQAKDIMDAGKLVT
##          1          .          .          .          .          .          60
##
##          61          .          .          .          .          .          120
## pdb|1AKE|A  DELVIALVKERIAQEDCRNGFLLDGFPRTPQADAMKEAGINVDYVLEFDVPDELIVDRI
##          61          .          .          .          .          .          120
##
##          121         .          .          .          .          .          180
## pdb|1AKE|A  VGRRVHAPSGRVYHVKNPPKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQM TAPLIG
##          121         .          .          .          .          .          180
##
##          181         .          .          .          214
## pdb|1AKE|A  YYSKEAEAGNTKYAKVDGTPVAEVRADLEKILG
##          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?

There are 218 amino acids are in this sequence.

BLAST search the PDB to find similar sequences and structures.

```

# Blast or hmmer search
b <- blast.pdb(aa)

```

```

## Searching ... please wait (updates every 5 seconds) RID = SBDPUY2X013
## ....
## Reporting 100 hits

```

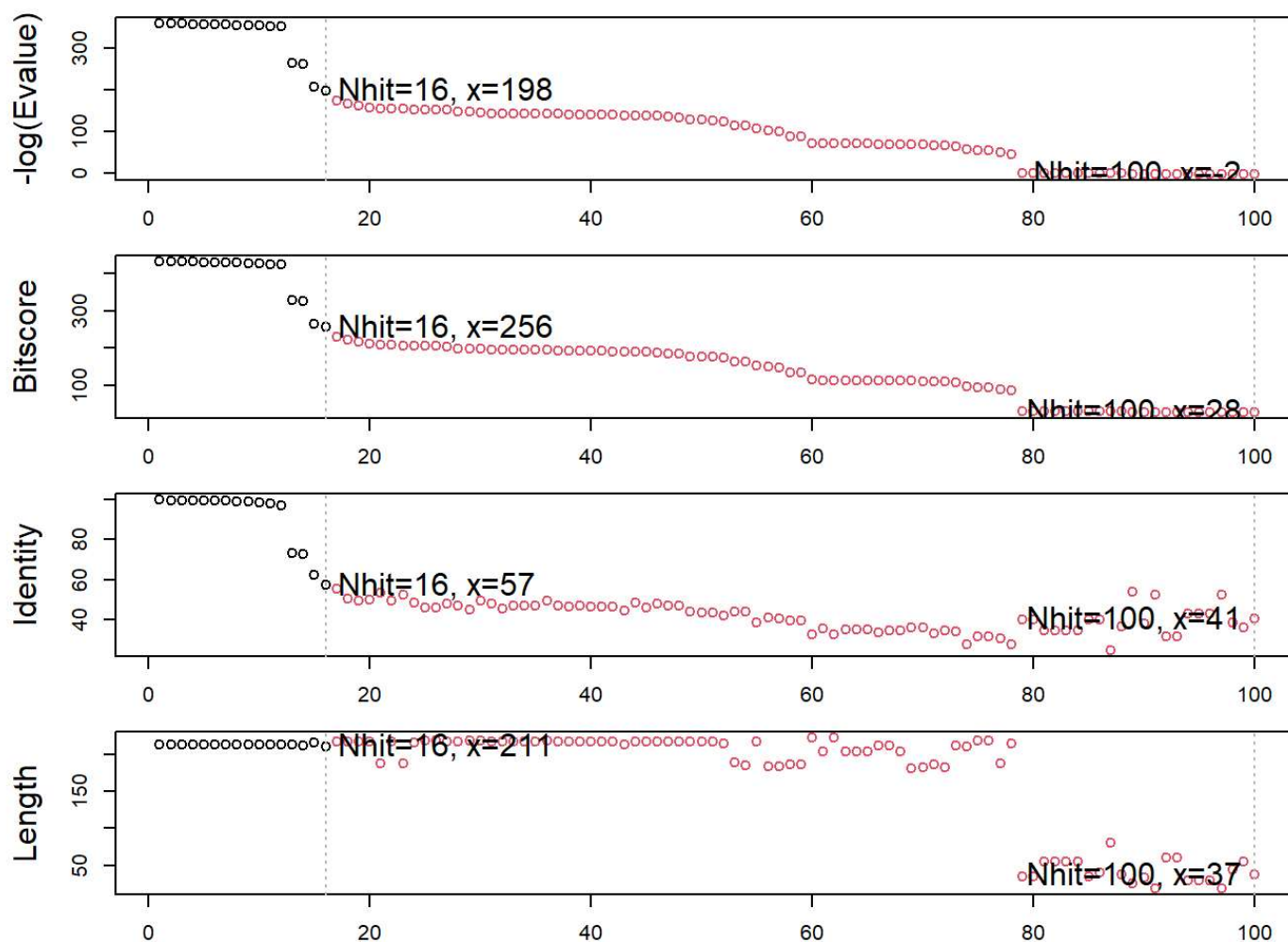
```

# Plot a summary of search results
hits <- plot(b)

```



```
## * Possible cutoff values: 197 -3
##      Yielding Nhits: 16 100
##
## * Chosen cutoff value of: 197
##      Yielding Nhits: 16
```



```
# List out some 'top hits'
hits$pdb.id <- c('1AKE_A', '4X8M_A', '6S36_A', '6RZE_A', '4X8H_A', '3HPR_A', '1E4V_A', '5EJE_A', '1E4Y_A', '3X2S_A', '6HAP_A', '6HAM_A', '4K46_A', '4NP6_A', '3GMT_A', '4PZL_A')
```

Identified a number of related PDB sequences to the query sequence.

```
# Download related PDB files
# get.pdb() and pdbslit() to fetch and parse the identified 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/  
## 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
```

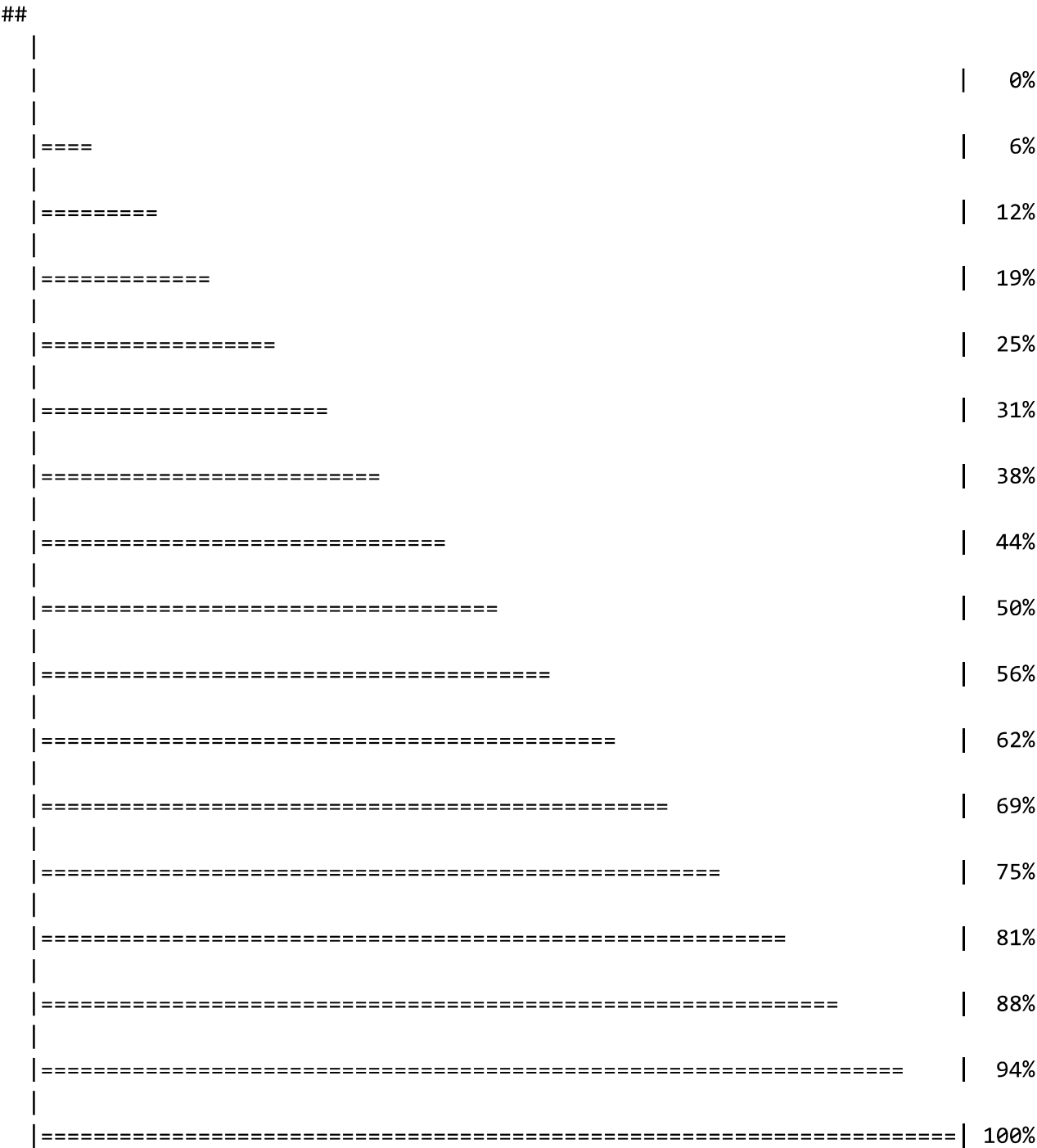
```
## Warning in get.pdb(hits$ pdb.id, path = "pdbs", split = TRUE, gzip = TRUE): pdbs/  
## 6HAM.pdb exists. Skipping download
```

```
## 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
```



Align and superpose structures

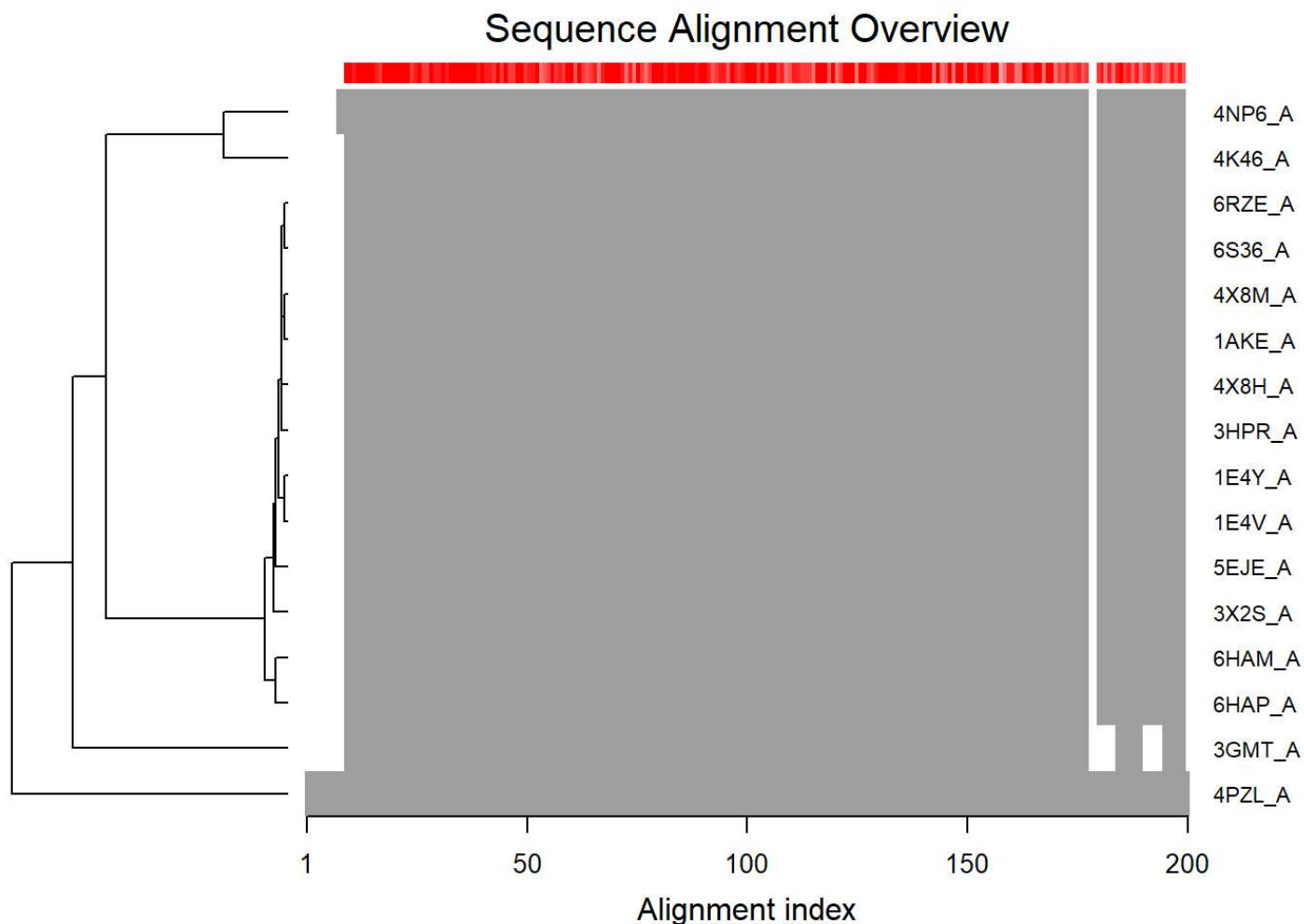
```
# Align related PDBs  
pdbs <- pdbaln(files, fit = TRUE, exefile="msa")
```

```
## 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
```

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

# Draw schematic alignment
plot(pdb, labels=ids)
```



The figure is displaying the schematic representation of alignment. Grey regions depict aligned residues, white regions depict gap regions. The red bar at the top depicts sequence conservation.

Optional: Viewing our superposed structures We can view the superposed results with `bio3d.view` `view()` function

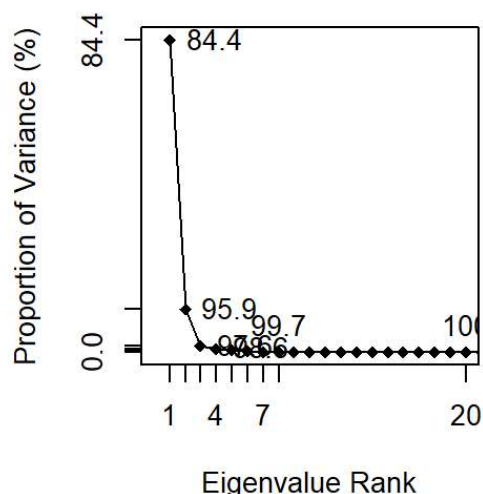
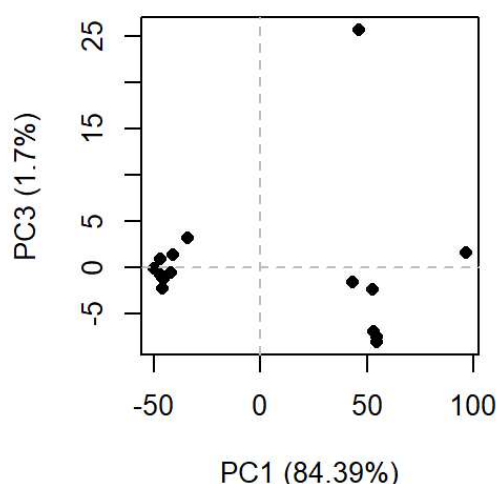
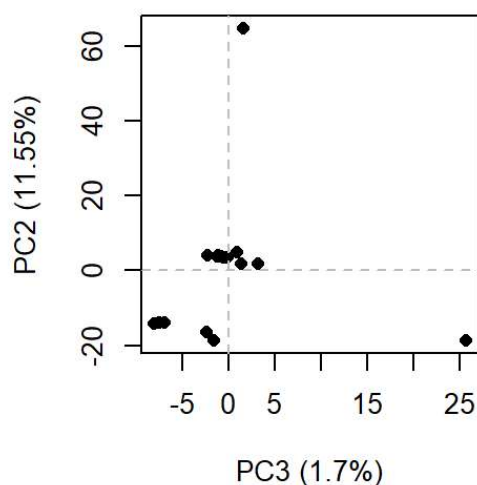
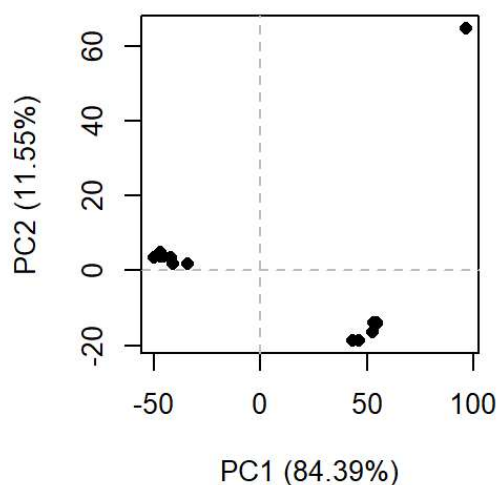
```
library(bio3d.view)
library(rgl)

view.pdb(pdb)
```

A 3D view of superposed ADK structures popped up.

Principal component analysis

```
# Perform PCA
pc.xray <- pca(pdbbs)
plot(pc.xray)
```



PCA results on Adenylate kinase X-ray structures. Each dot represents one PDB structure

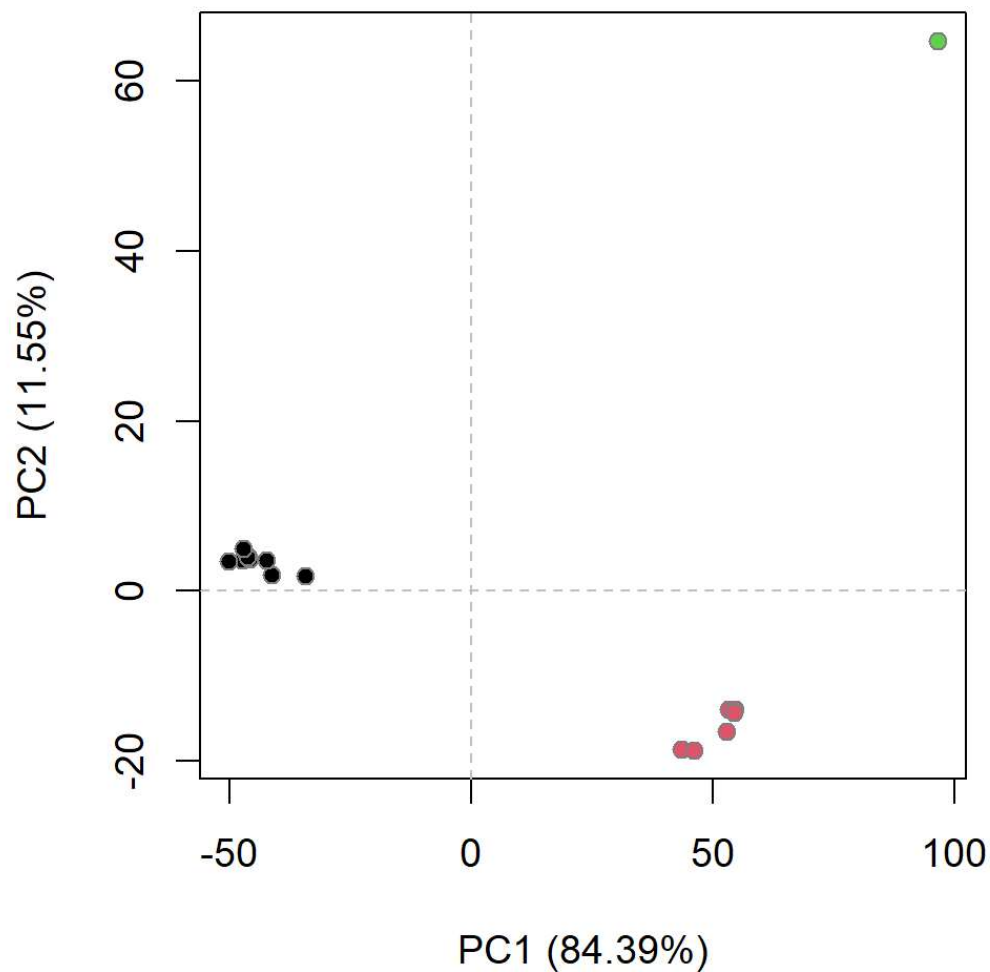
rmsd() will calculate all pairwise RMSD values of the structural ensemble

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

```
## Warning in rmsd(pdbbs): 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)
```



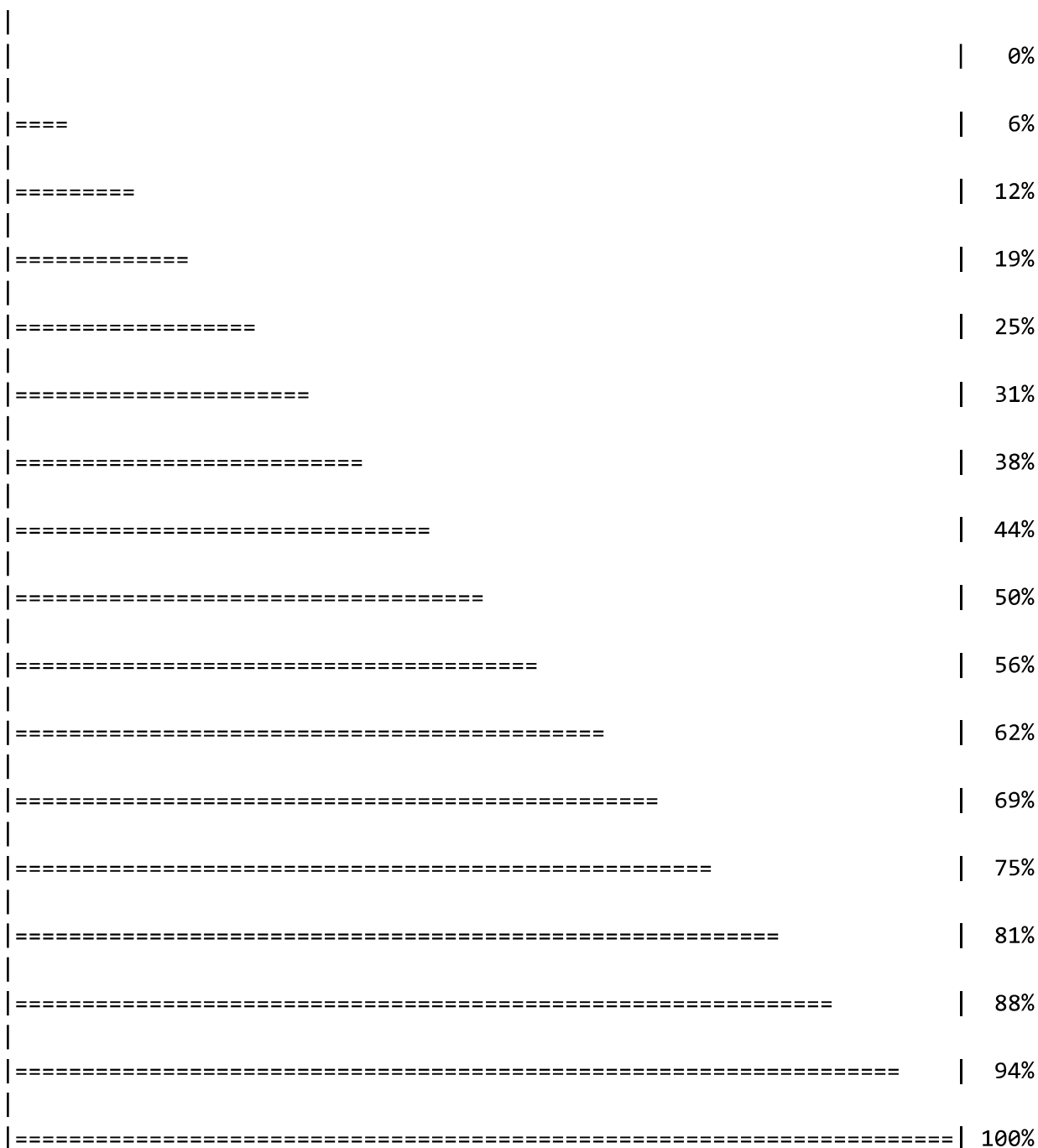
Normal mode analysis

nma() provides normal mode analysis which facilitates characterizing and comparing flexibility profiles of related protein structures.

```
# NMA of all structures
modes <- nma(pdbbs)
```

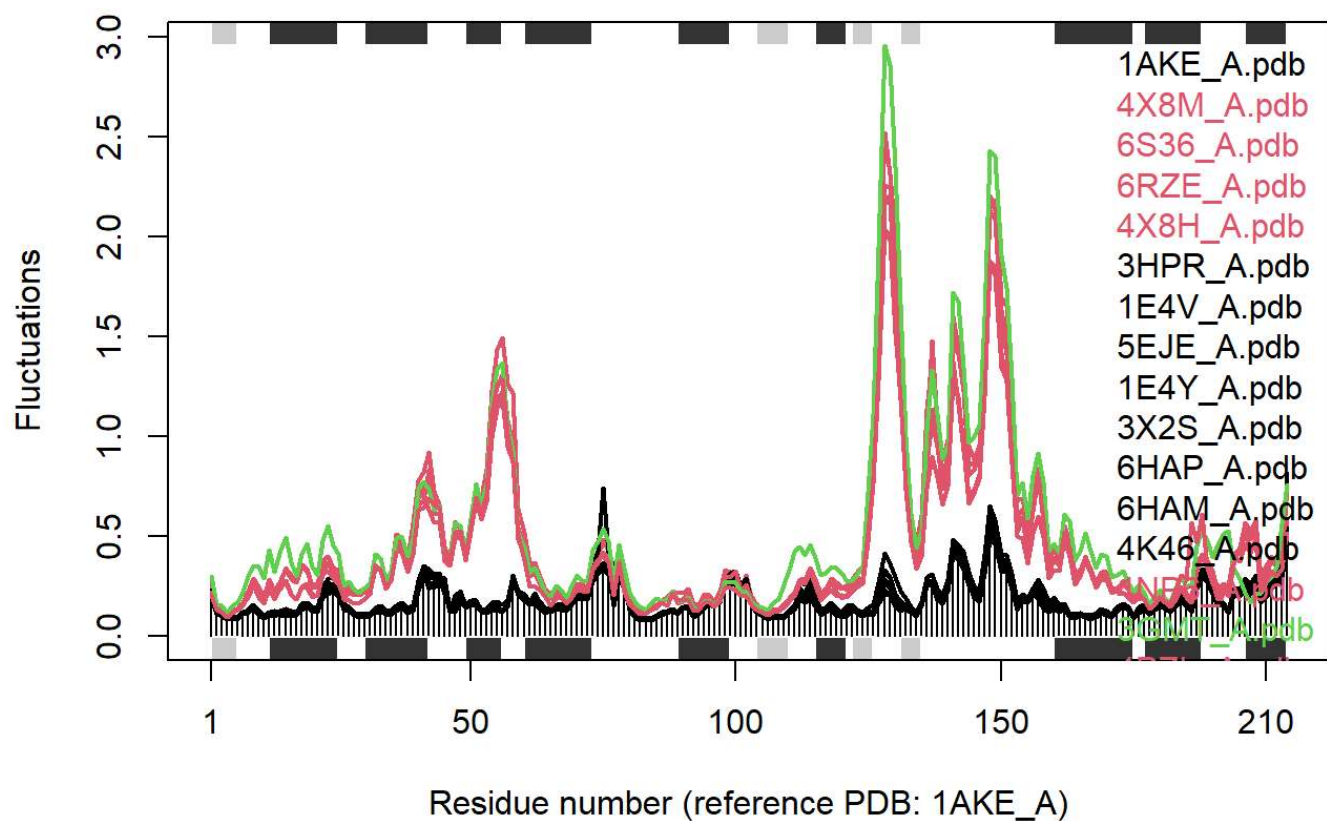


```
##
## Details of Scheduled Calculation:
## ... 16 input structures
## ... storing 606 eigenvectors for each structure
## ... dimension of x$U.subspace: ( 612x606x16 )
## ... coordinate superposition prior to NM calculation
## ... aligned eigenvectors (gap containing positions removed)
## ... estimated memory usage of final 'eNMA' object: 45.4 Mb
##
##
```



```
plot(modes, pdirs, col=grps.rd)
```

```
## Extracting SSE from pdb$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?

This x axis probably demonstrating the differences in the sequence, in nucleotide binding or amino acids. The black and colored lines are different. They differ most in residue numbers 125 to 150