

Class 11 Structural Bioinformatics (Pt. 1)

Cindy Tran

1/22/2022

1. Introduction to the RCSB Protein Data Bank

PDB Statistics

Download a CSV file from the PDB site (accessible from “Analyze” > “PDB Statistics” > “by Experimental Method and Molecular Type”. Move this CSV file into your RStudio project and use it to answer the following questions:

```
read.csv("Data Export Summary.csv")
```

##	Molecular.Type	X.ray	NMR	EM	Multiple.methods	Neutron	Other	
## 1	Protein (only)	143950	11863	6571		179	70	32
## 2	Protein/Oligosaccharide	8514	31	1086		5	0	0
## 3	Protein/NA	7610	274	2127		3	0	0
## 4	Nucleic acid (only)	2393	1396	61		8	2	1
## 5	Other	150	31	3		0	0	0
## 6	Oligosaccharide (only)	11	6	0		1	0	4
##	Total							
## 1	162665							
## 2	9636							
## 3	10014							
## 4	3861							
## 5	184							
## 6	22							

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

```
(143950 + 8514 + 7610 + 2393 + 150 + 11) / (162665 + 9636 + 10014 + 3861 + 184 + 22) * 100
```

```
## [1] 87.25521
```

```
(6571 + 1086 + 2127 + 61 + 3) / (162665 + 9636 + 10014 + 3861 + 184 + 22) * 100
```

```
## [1] 5.283772
```

About 87% by X-ray and 5% by electron microscopy (about 92% by either X-ray or electron microscopy)

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

```
162665 / (162665 + 9636 + 10014 + 3861 + 184 + 22) * 100
```

```
## [1] 87.27506
```

About 87% are protein only.

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?

1856

2. Visualizing the HIV-1 Protease Structure

Using Atom Selections

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

It is because we chose to represent the entire water molecule as one sphere.

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)?

127

Sequence Viewer Extension [OPTIONAL]

Q6. As you have hopefully observed HIV protease is a homodimer (i.e. it is composed of two identical chains). With the aid of the graphic display and the sequence viewer extension can you identify secondary structure elements that are likely to only form in the dimer rather than the monomer?

Coil

3. Introduction to Bio3D in R

```
library(bio3d)
```

Reading PDB File Data into R

```
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"
```

```
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>
```

4. Comparative Structure Analysis of Adenylate Kinase

Setup

```
# 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

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

```
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  MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMLRAAVKSGSELGKQAKDIMDAGKLV
##          1          .          .          .          .          .          60
##
##          61          .          .          .          .          .          120
## pdb|1AKE|A  DELVIALVKERIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFDVPDELIVDRI
##          61          .          .          .          .          .          120
##
##          121         .          .          .          .          .          180
## pdb|1AKE|A  VGRRVHAPSGRVYHVKFNPPKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQMTAPLIG
##          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?

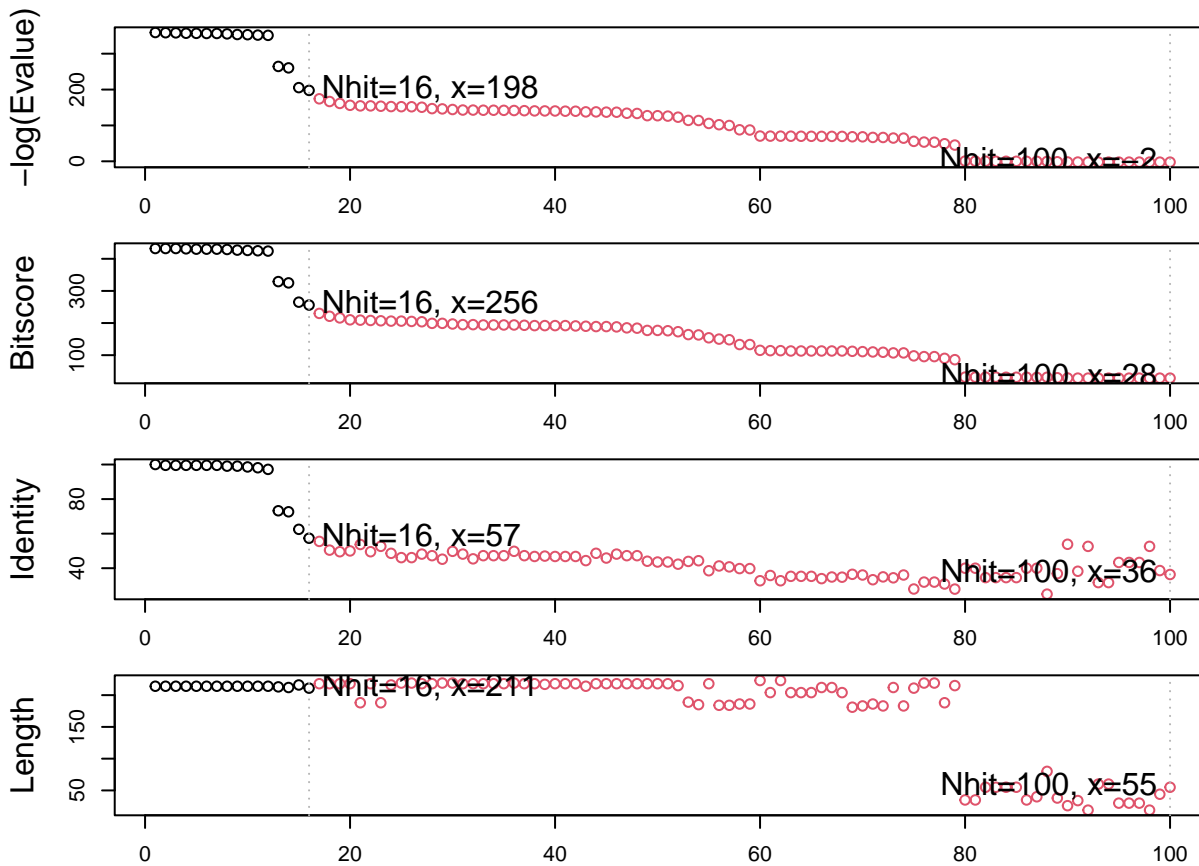
214

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

```
## Searching ... please wait (updates every 5 seconds) RID = ONWCTYT6013
## ....
## 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'
head(hits$ pdb.id)
```

```
## [1] "1AKE_A" "4X8M_A" "6S36_A" "6RZE_A" "4X8H_A" "3HPR_A"
```

```
# Download related PDB files
```

```
files <- get.pdb(hits$ pdb.id, path="pdb", split=TRUE, gzip=TRUE)
```

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

```
## Warning in get.pdb(hits$ pdb.id, path = "pdb", split = TRUE, gzip = TRUE): pdb/
## 4X8M.pdb exists. Skipping download
```

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

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

```
## Warning in get.pdb(hits$ pdb.id, path = "pdb", split = TRUE, gzip = TRUE): pdb/
## 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
```

```

## pdb/split_chain/1E4Y_A.pdb
## pdb/split_chain/3X2S_A.pdb
## pdb/split_chain/6HAP_A.pdb
## pdb/split_chain/6HAM_A.pdb
## pdb/split_chain/4K46_A.pdb
## pdb/split_chain/4NP6_A.pdb
## pdb/split_chain/3GMT_A.pdb
## pdb/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: pdb/split_chain/1AKE_A.pdb
##   PDB has ALT records, taking A only, rm.alt=TRUE
## pdb/seq: 2   name: pdb/split_chain/4X8M_A.pdb
## pdb/seq: 3   name: pdb/split_chain/6S36_A.pdb
##   PDB has ALT records, taking A only, rm.alt=TRUE
## pdb/seq: 4   name: pdb/split_chain/6RZE_A.pdb
##   PDB has ALT records, taking A only, rm.alt=TRUE
## pdb/seq: 5   name: pdb/split_chain/4X8H_A.pdb
## pdb/seq: 6   name: pdb/split_chain/3HPR_A.pdb
##   PDB has ALT records, taking A only, rm.alt=TRUE
## pdb/seq: 7   name: pdb/split_chain/1E4V_A.pdb
## pdb/seq: 8   name: pdb/split_chain/5EJE_A.pdb
##   PDB has ALT records, taking A only, rm.alt=TRUE
## pdb/seq: 9   name: pdb/split_chain/1E4Y_A.pdb
## pdb/seq: 10  name: pdb/split_chain/3X2S_A.pdb
## pdb/seq: 11  name: pdb/split_chain/6HAP_A.pdb
## pdb/seq: 12  name: pdb/split_chain/6HAM_A.pdb
##   PDB has ALT records, taking A only, rm.alt=TRUE
## pdb/seq: 13  name: pdb/split_chain/4K46_A.pdb
##   PDB has ALT records, taking A only, rm.alt=TRUE
## pdb/seq: 14  name: pdb/split_chain/4NP6_A.pdb
## pdb/seq: 15  name: pdb/split_chain/3GMT_A.pdb
## pdb/seq: 16  name: pdb/split_chain/4PZL_A.pdb

```

```

# Vector containing PDB codes for figure axis

```

```

ids <- basename.pdb(pdb$id)

```

```

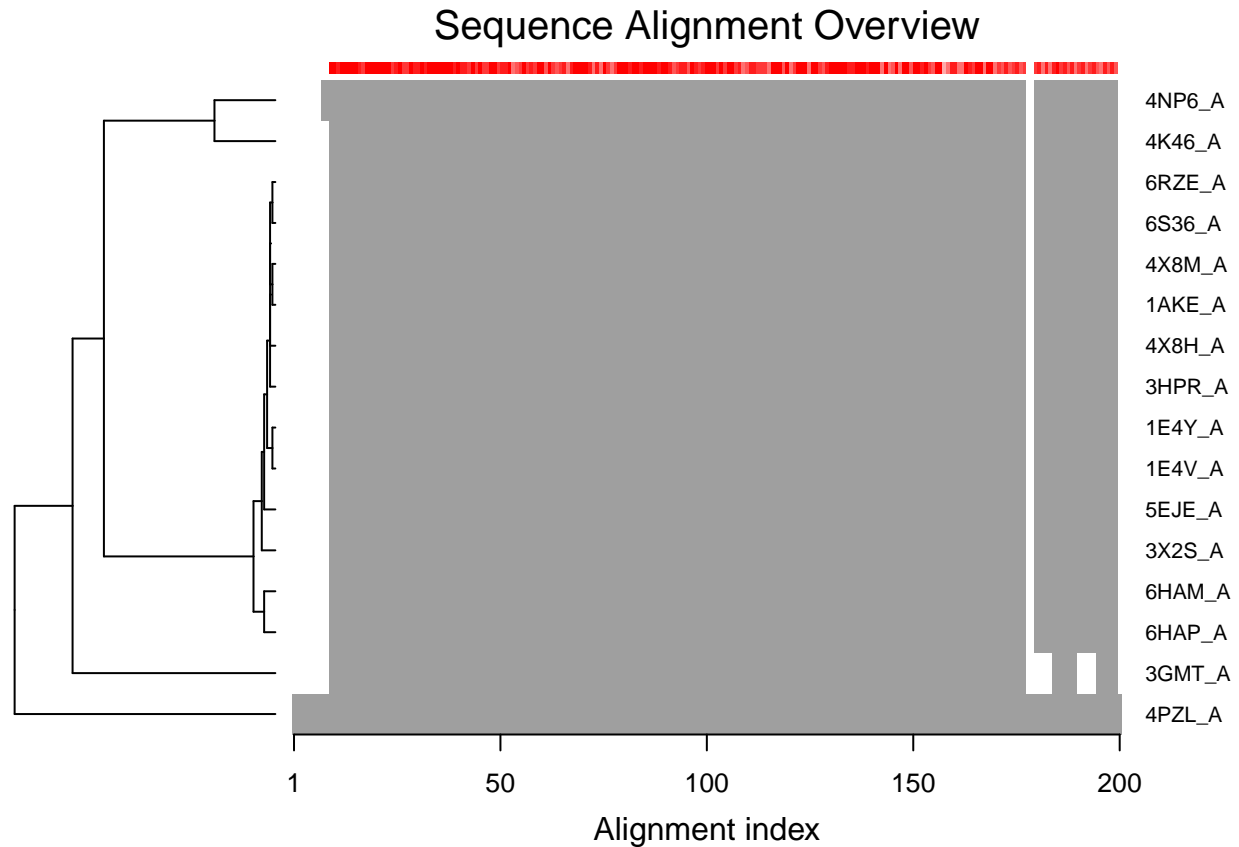
# Draw schematic alignment

```

```

plot(pdb, labels = ids)

```

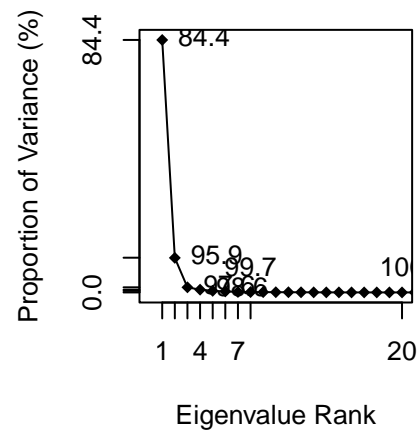
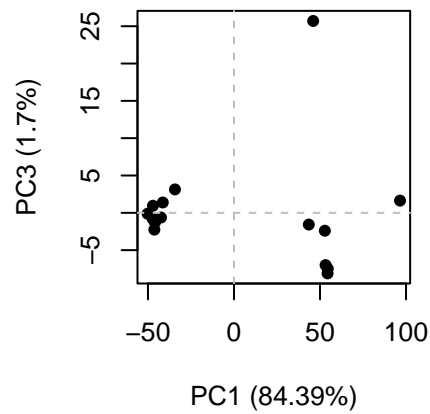
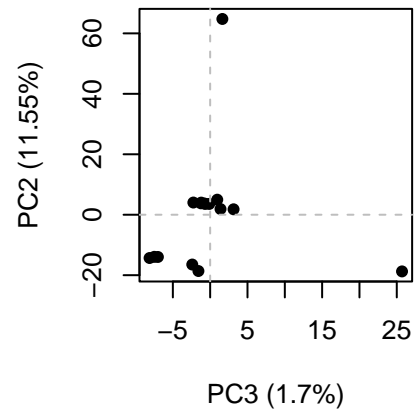
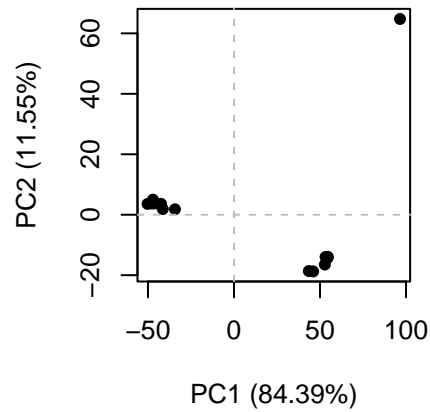
Optional: Viewing our superposed structures

```
library(bio3d.view)
#install.packages("rgl")
library(rgl)

view.pdbs(pdbs)
```

Principal Component Analysis

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

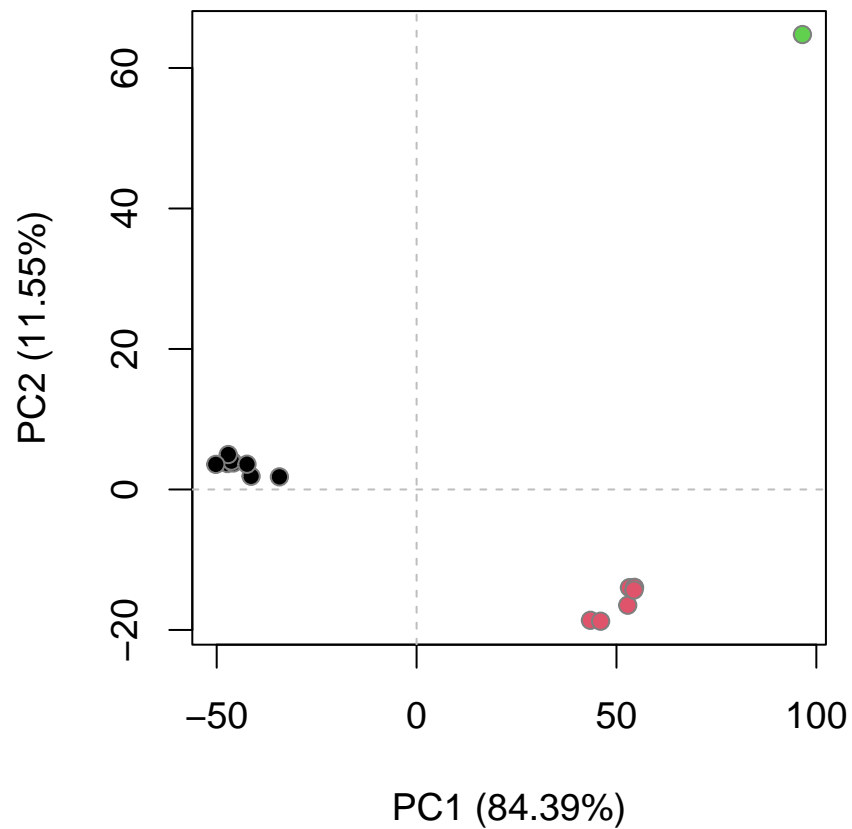


```
# Calculate RMSD
rd <- rmsd(pdb)
```

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



5. Optional Further Visualization

```
# Visualize first principal component
pc1 <- mktrj(pc.xray, pc=1, file="pc_1.pdb")
```

```
view.xyz(pc1)
```

```
## Potential all C-alpha atom structure(s) detected: Using calpha.connectivity()
```

```
view.xyz(pc1, col=vec2color( rmsf(pc1) ))
```

```
## Potential all C-alpha atom structure(s) detected: Using calpha.connectivity()
```

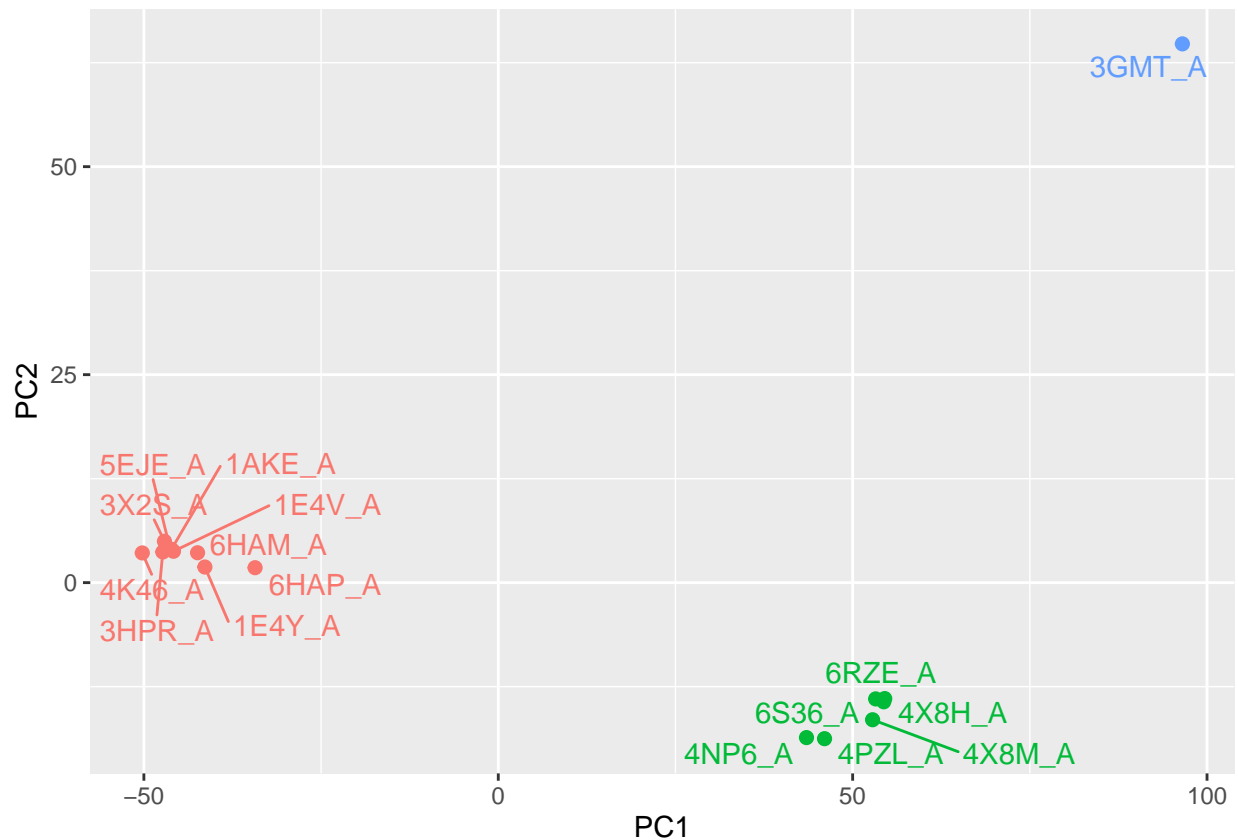
```
#Plotting results with ggplot2
```

```
library(ggplot2)
```

```
library(ggrepel)
```

```
df <- data.frame(PC1=pc.xray$z[,1],
                 PC2=pc.xray$z[,2],
                 col=as.factor(grps.rd),
                 ids=ids)
```

```
p <- ggplot(df) +
  aes(PC1, PC2, col=col, label=ids) +
  geom_point(size=2) +
  geom_text_repel(max.overlaps = 20) +
  theme(legend.position = "none")
p
```



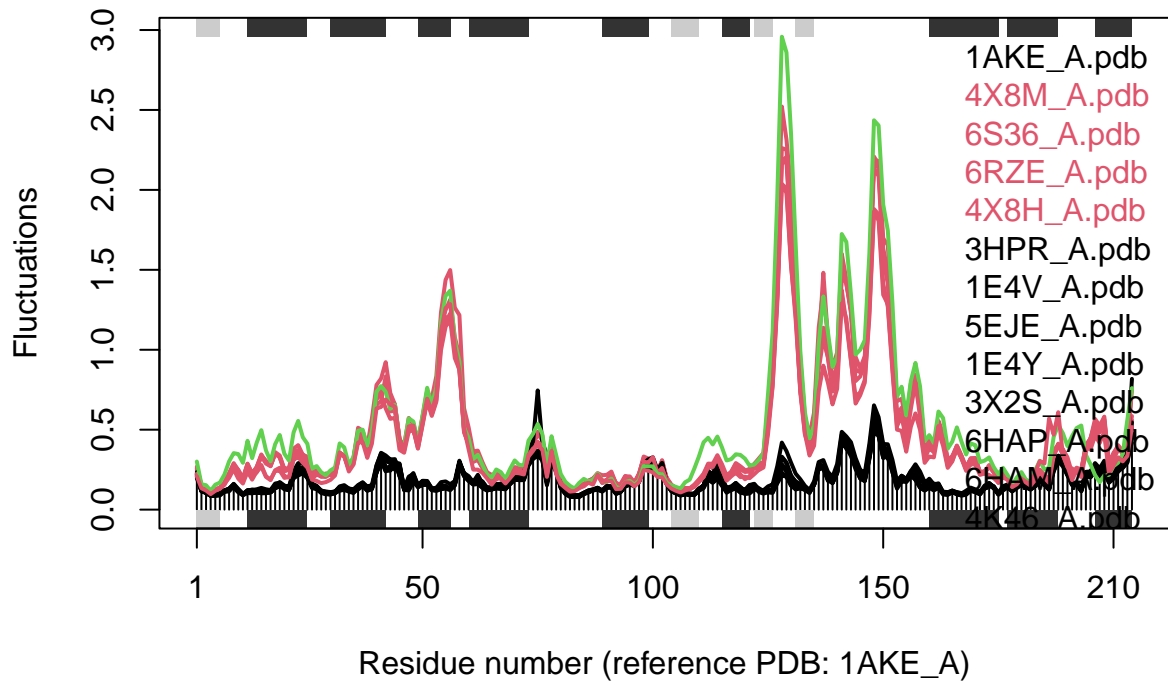
6. Normal Mode Analysis

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

```
##
## 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, 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 are similar in terms of overall shape, but they differ in that the black lines have less fluctuations overall than the colored lines. I think they differ most along residues 40-60 and residues 125-155. I think this is because these regions of residues are where the protein is most flexible and can change conformations.