

Class 9: Structural Bioinformatics

Elena

Table of contents

Introduction to the RCSB Protein Data Bank (PDB)	2
Q1. What percentage of structures in the PDB are solved by X-Ray and Electron Microscopy.	2
Q2. What proportion of structures in the PDB are protein?	3
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?	3
Visualizing the HIV-1 protease structure	4
Q4. Water molecules normally have 3 atoms. Why do we see just one atom per water molecule in this structure?	5
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	5
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.	7
Introduction to Bio3D in R	7
Q7. How many amino acid residues are there in this pdb object?	8
Q8. Name one of the two non-protein residues?	8
Q9. How many protein chains are in this structure?	8
Predicting functional motions of a single structure	9
Comparative structure analysis of Adenylate Kinase	10
Q10. Which of the packages above is found only on BioConductor and not CRAN? .	11
Q11. Which of the above packages is not found on BioConductor or CRAN?	11
Q12. True or False? Functions from the devtools package can be used to install packages from GitHub and BitBucket?	11

Search and retrieve ADK structures	11
Q13. How many amino acids are in this sequence, i.e. how long is this sequence? . .	12
Align and superpose structures	15
PCA	17
Normal mode analysis	18
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?	20

Introduction to the RCSB Protein Data Bank (PDB)

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

	Molecular.Type	X.ray	NMR	EM	Multiple.methods	Neutron	Other
1	Protein (only)	150,417	12,056	8,586	188	72	32
2	Protein/Oligosaccharide	8,869	32	1,552	6	0	0
3	Protein/NA	7,943	280	2,690	6	0	0
4	Nucleic acid (only)	2,522	1,425	74	13	2	1
5	Other	154	31	6	0	0	0
6	Oligosaccharide (only)	11	6	0	1	0	4
	Total						
1		171,351					
2		10,459					
3		10,919					
4		4,037					
5		191					
6		22					

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

92.8%

```
#Need to get rid of commas in numbers
csv_X.ray <- gsub(",", "", csv$X.ray)
csv_EM <- gsub(",", "", csv$EM)
csv_Total <- gsub(",", "", csv$Total)
```

```
(sum(as.numeric(csv_X.ray))+sum(as.numeric(csv_EM)))/sum(as.numeric(csv_Total))
```

```
[1] 0.9281395
```

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

87.0%

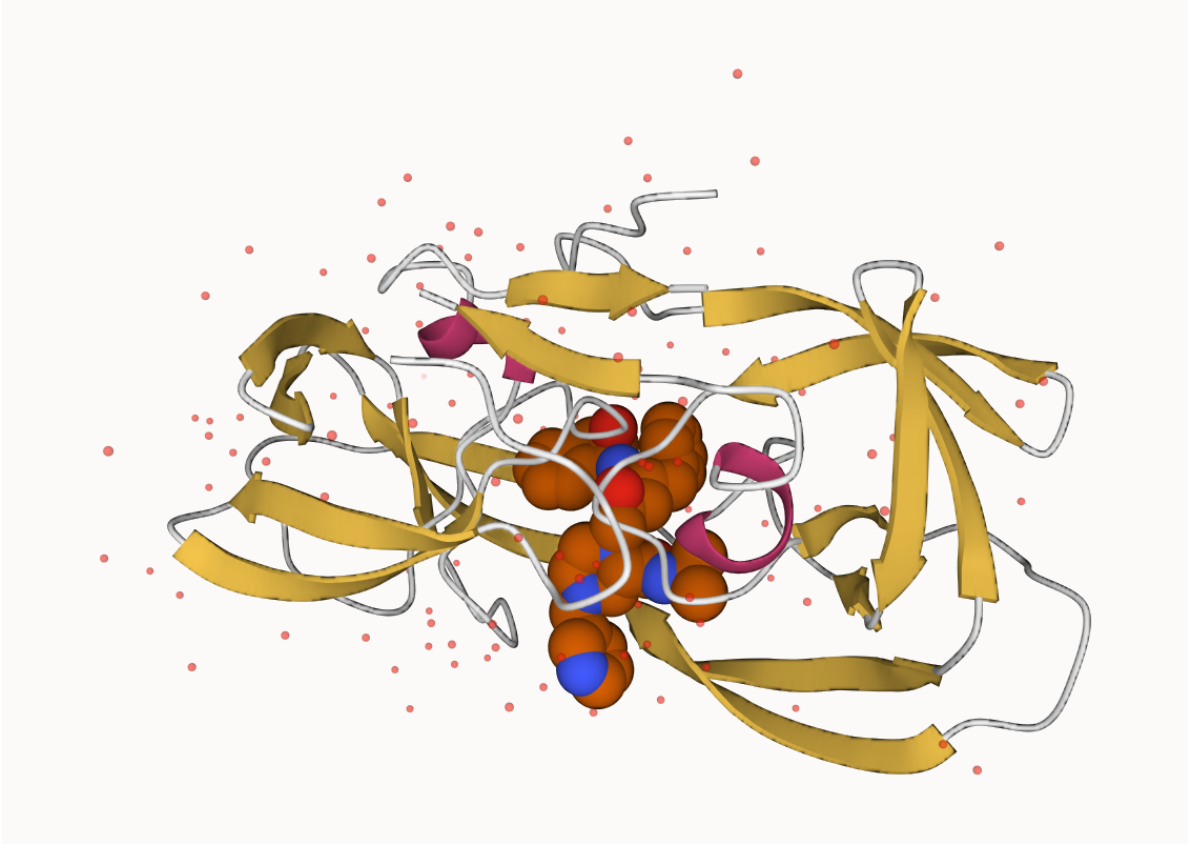
```
as.numeric(csv_Total[1])/sum(as.numeric(csv_Total))
```

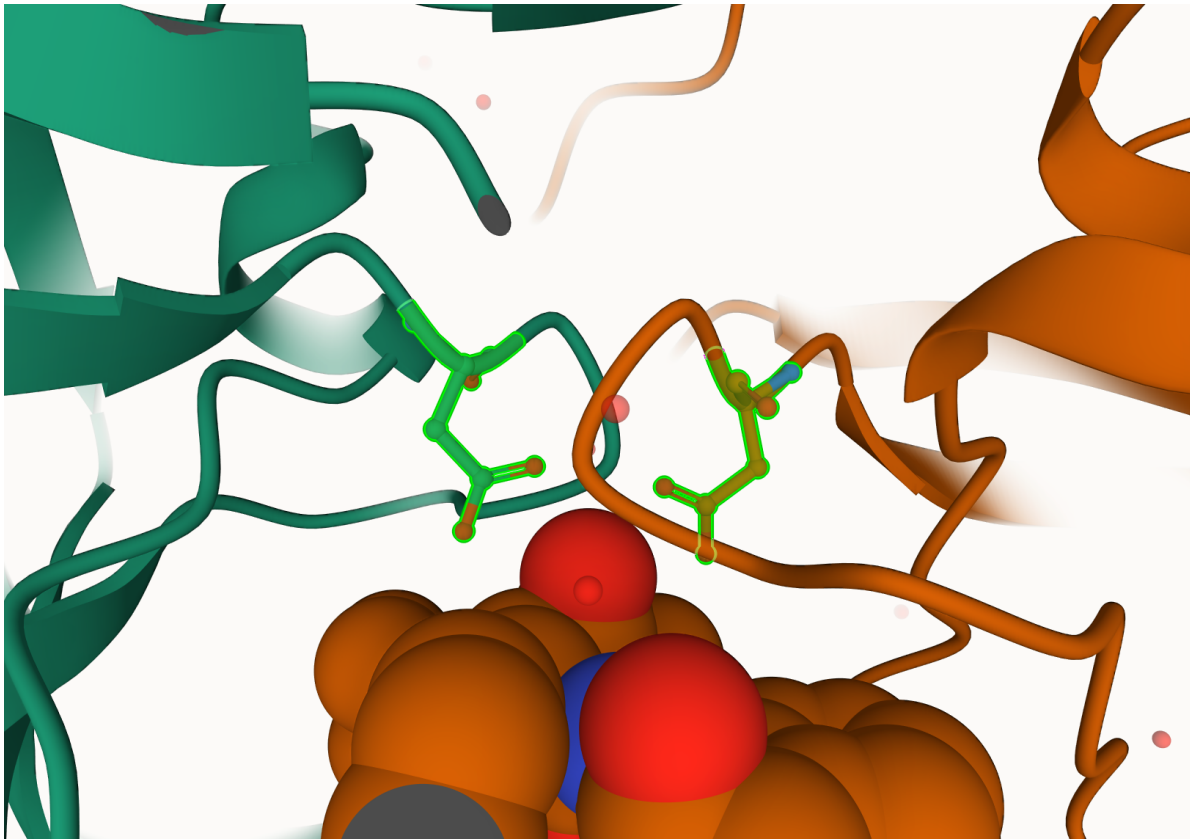
```
[1] 0.8698948
```

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?

Searched for HIV AND Protease, filtered for proteins, and yielded 1264 results. Each structure seems to be bound to different ligands.

Visualizing the HIV-1 protease structure





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

Only see oxygen atom because hydrogen atoms are much smaller than oxygen.

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

Residue 308.

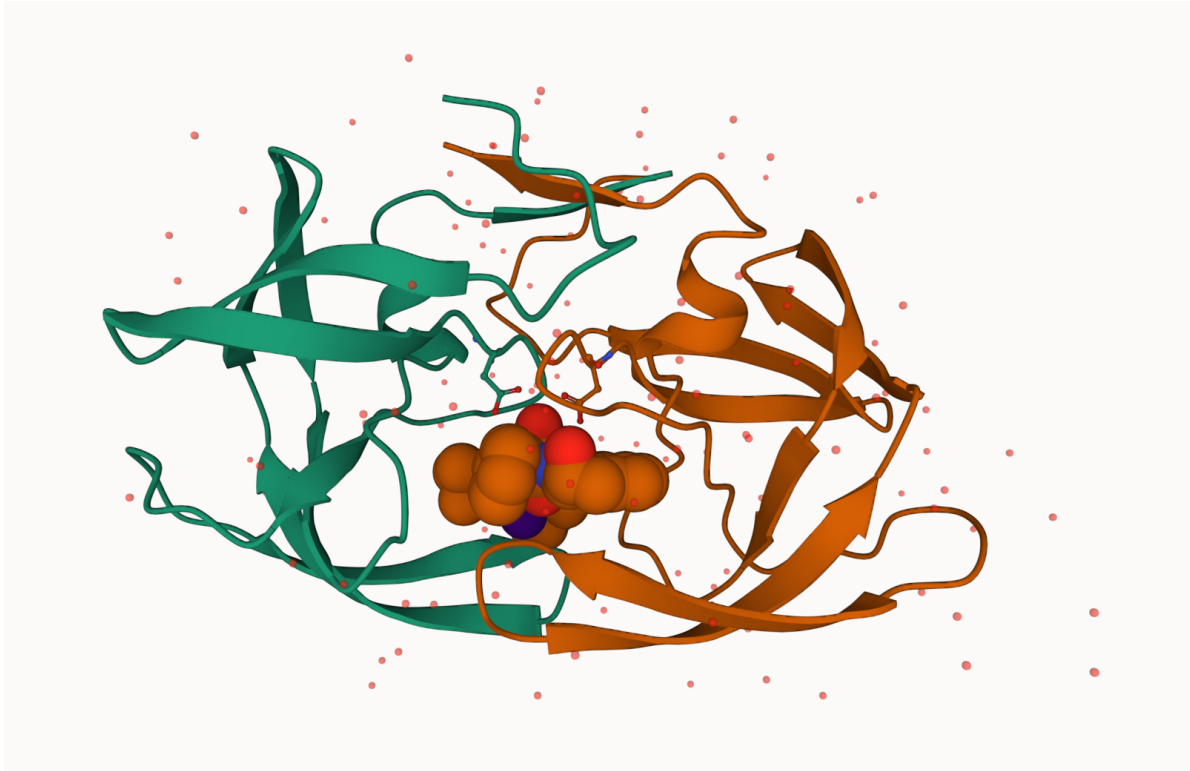


Figure 1: HIV protease with ligand (spacefill), catalytic residues (ball and stick), and critical H₂O molecule (space fill, dark purple)

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.

Introduction to Bio3D in 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
```

```
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>

Q7. How many amino acid residues are there in this pdb object?

198.

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

HOH (H₂O).

Q9. How many protein chains are in this structure?

2 chains.

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:

```
MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMRLRAAVKSGSELGKQAKDIMDAGKLV
TDELVIALVKERIAQEDCRNGFLDGFPRTPQADAMKEAGINVDYVLEFDVPDELIVDKI
VGRRVHAPSGRVYHVKFNPVKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQMTAPLIG
YYSKEAEAGNTKYAKVDGTPVAEVRADLEKILG
```

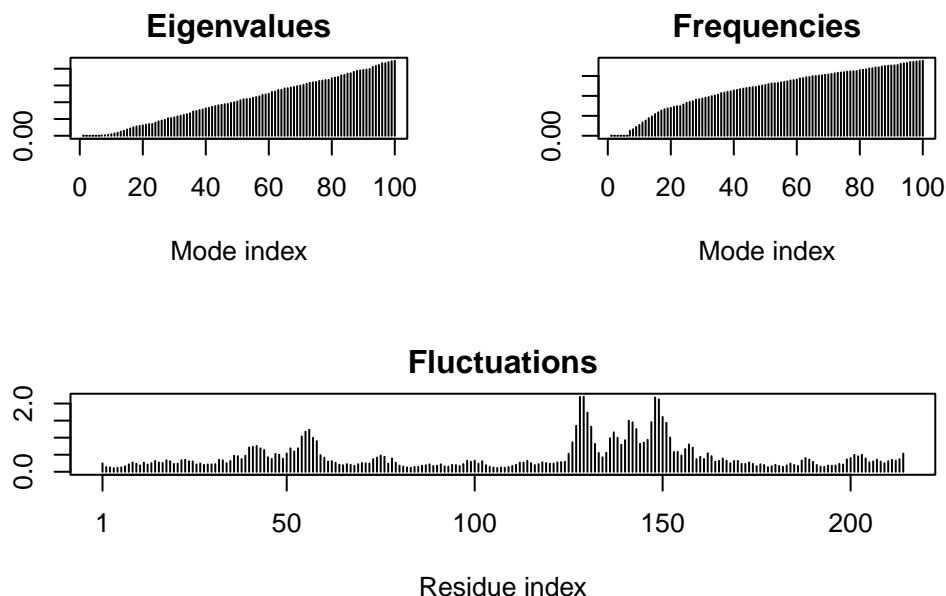
```
+ attr: atom, xyz, seqres, helix, sheet,
      calpha, remark, call
```

NMA to predict protein flexibility and potential functional motions (conformational changes)

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

```
Building Hessian...      Done in 0.086 seconds.
Diagonalizing Hessian... Done in 0.264 seconds.
```

```
plot(m)
```



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

Comparative structure analysis of Adenylate Kinase

“The `install.packages()` function is used to install packages from the main CRAN repository for R packages. BioConductor is a separate repository of R packages focused on high-throughput biomolecular assays and biomolecular data. Packages from BioConductor can be installed using the `BiocManager::install()` function. Finally, R packages found on GitHub or BitBucket can be installed using `devtools::install_github()` and `devtools::install_bitbucket()` functions.”

```
#install.packages("bio3d")
#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("lake_A")
```

Warning in get.seq("lake_A"): Removing existing file: seqs.fasta

Fetching... Please wait. Done.

```
aa
```

```

      1      .      .      .      .      .      60
pdb|1AKE|A  MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMRLAAVKSSELGKQAKDIMDAGKLV
      1      .      .      .      .      .      60

      61      .      .      .      .      .      120
pdb|1AKE|A  DELVIALVKERIAQEDCRNGFLLDGFRTIPQADAMKEAGINVDYVLEFDVPDELIVDRI
      61      .      .      .      .      .      120

     121      .      .      .      .      .      180
pdb|1AKE|A  VGRRVHAPSGRVYHVKNPPKVEGKDDVTGEELTRKDDQEETVRKRLVEYHQMTPALIG
     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 hmmer search
b <- blast.pdb(aa)

```

```

Searching ... please wait (updates every 5 seconds) RID = NZ0BU9UD013
.....
Reporting 98 hits

```

```

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

```

```

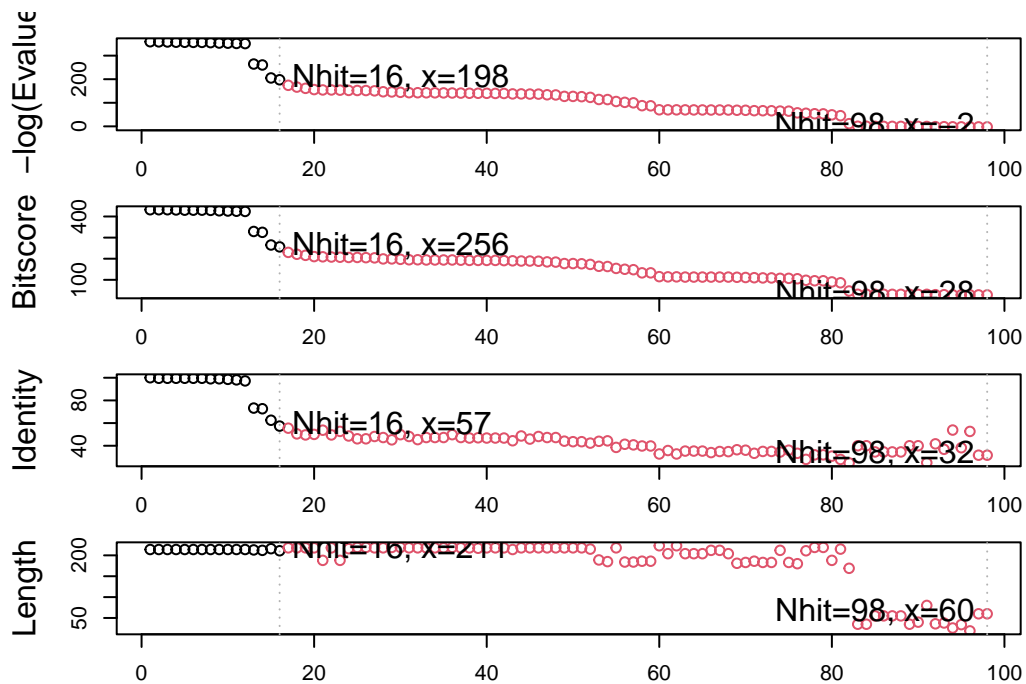
* Possible cutoff values: 197 -3
  Yielding Nhits:       16 98

```

```

* 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"
```

```
#There are 16 hits, but will use the ones listed on the workbook
```

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

```
# 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.gz exists. Skipping download
```

```
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%



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/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

```

```

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

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

#Worked!

```

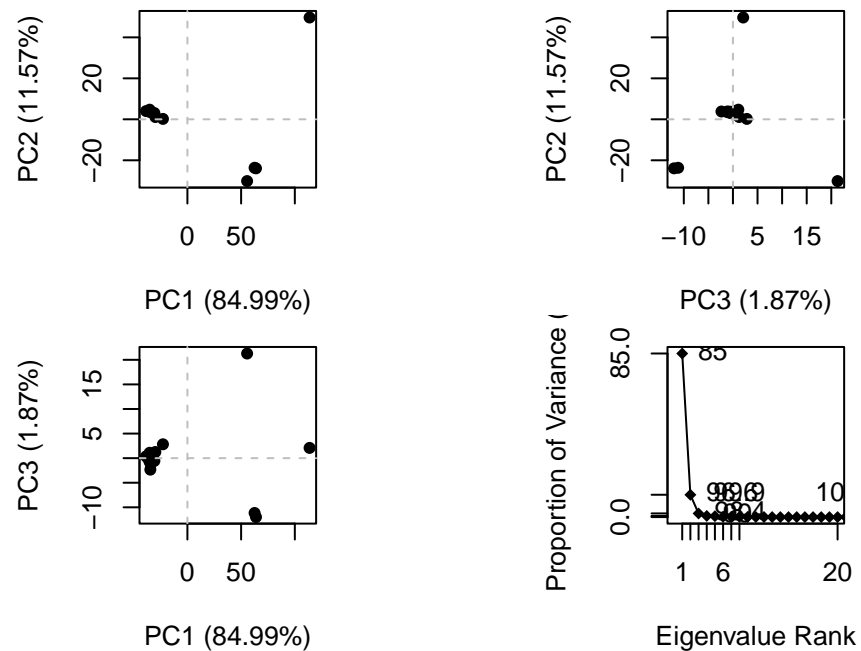


```
#But is causing issues with PDF export so not including in the report
```

```
#anno <- pdb.annotate(ids)  
#unique(anno$source)
```

PCA

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



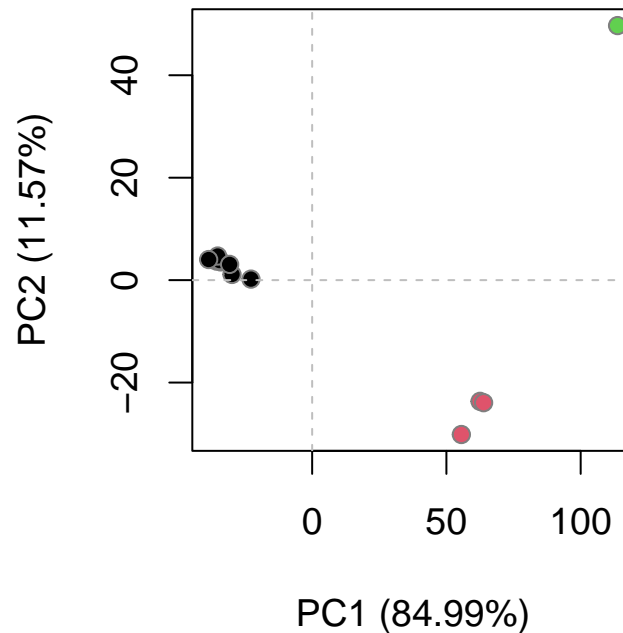
```
#Visualize first principal component  
#mktrj(pc.xray, pc=1, file="pc_1.pdb")
```

```
#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)
```

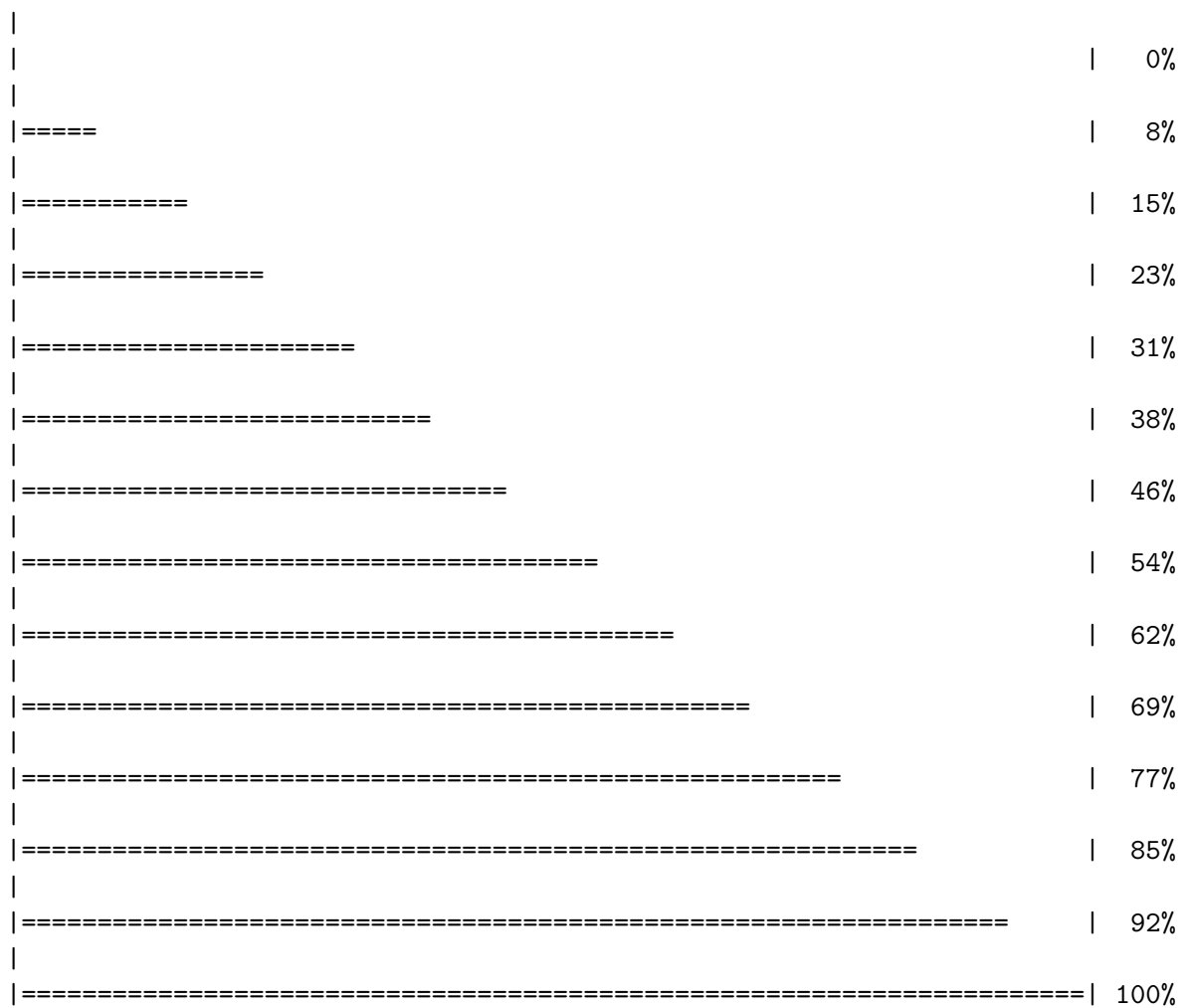


Normal mode analysis

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

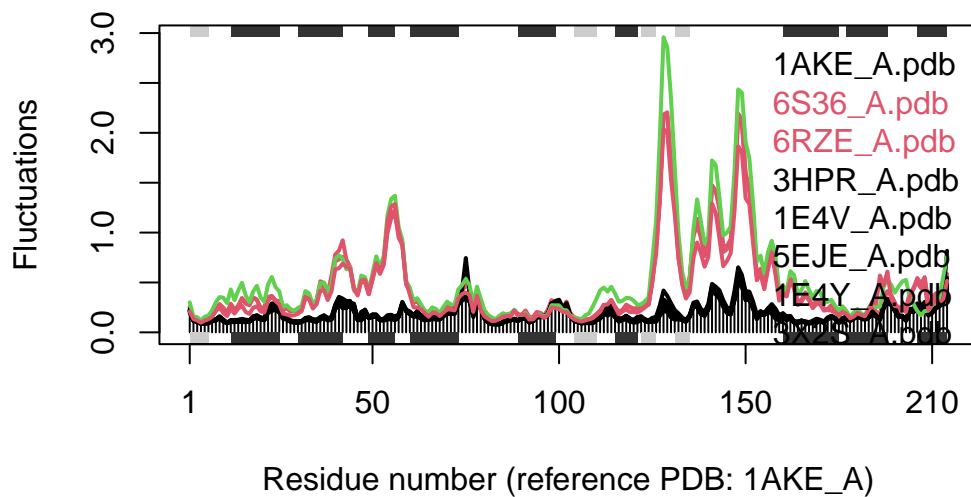
Details of Scheduled Calculation:

```
... 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
```



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

Extracting SSE from pdbc\$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?

There are regions where the curves are similar and regions that are different. Perhaps different most in sites that are not the active site. These proteins are similar supposedly because they share similar catalytic activity as conferred by the active site.