

Class 10: Structural Bioinformatics 1

Rachel Lamm (A18518313)

Table of contents

PDB Statistics	1
Visualizing the HIV-1 protease structure	4
Bio3D package for structural bioinformatics	6
Quick PDB visualization in R	8
Predicting functional motions of a single structure	8
Comparative structure analysis of Adenylate Kinase	11
Align and superpose structures	12
Principal component analysis	18
Normal mode analysis [optional]	21

PDB Statistics

The Protein Data Bank (PDB) is the main repository of biomolecular structures. Let's see what it contains"

```
stats<- read.csv("pdbstats26")
stats
```

	Molecular.Type	X.ray	EM	NMR	Integrative	Multiple.methods
1	Protein (only)	178,795	21,825	12,773	343	226
2	Protein/Oligosaccharide	10,363	3,564	34	8	11
3	Protein/NA	9,106	6,335	287	24	7
4	Nucleic acid (only)	3,132	221	1,566	3	15
5	Other	175	25	33	4	0
6	Oligosaccharide (only)	11	0	6	0	1
	Neutron	Other	Total			
1	84	32	214,078			
2	1	0	13,981			

```

3      0      0 15,759
4      3      1  4,941
5      0      0   237
6      0      4    22

```

```
stats$X.ray
```

```
[1] "178,795" "10,363" "9,106" "3,132" "175" "11"
```

```
sum(stats$Neutron)
```

```
[1] 88
```

The comma in these numbers leads to the numbers here being read as character.

```
library(readr)
stats<- read_csv("pdbstats26")
```

```
Rows: 6 Columns: 9
```

```
-- Column specification -----
```

```
Delimiter: ","
```

```
chr (1): Molecular Type
```

```
dbl (4): Integrative, Multiple methods, Neutron, Other
```

```
num (4): X-ray, EM, NMR, Total
```

```
i Use `spec()` to retrieve the full column specification for this data.
```

```
i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

```
stats
```

```
# A tibble: 6 x 9
```

```

`Molecular Type`  `X-ray`  EM  NMR Integrative `Multiple methods` Neutron
<chr>            <dbl> <dbl> <dbl>          <dbl>          <dbl> <dbl>
1 Protein (only)  178795 21825 12773          343          226   84
2 Protein/Oligosacch~ 10363 3564  34           8           11    1
3 Protein/NA      9106 6335  287          24           7     0
4 Nucleic acid (only) 3132  221 1566           3           15    3
5 Other           175   25  33            4           0     0
6 Oligosaccharide (o~  11    0   6            0           1     0

```

```
# i 2 more variables: Other <dbl>, Total <dbl>
```

```
sum(stats$`X-ray`)
```

```
[1] 201582
```

```
n.xray <-sum(stats$`X-ray`)
```

```
n.total <-sum(stats$Total)
```

```
n.xray/n.total
```

```
[1] 0.8095077
```

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

```
n.xray <- sum(stats$`X-ray`)
```

```
n.em <- sum(stats$EM)
```

```
n.total <- sum(stats$Total)
```

```
percent.xray <- (n.xray / n.total) * 100
```

```
percent.em <- (n.em / n.total) * 100
```

```
percent.xray
```

```
[1] 80.95077
```

```
percent.em
```

```
[1] 12.83843
```

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

```
n.protein <- sum(  
  stats$Total[stats$`Molecular Type` %in%  
    c("Protein (only)", "Protein/Oligosaccharide", "Protein/NA")]  
)
```

```
n.total <- sum(stats$Total)
```

```
n.protein / n.total
```

```
[1] 0.979118
```

Q3: Skip...Looking up HIV structures

Visualizing the HIV-1 protease structure

We can use the Molstar viewer online: <https://molstar.org/viewer/>

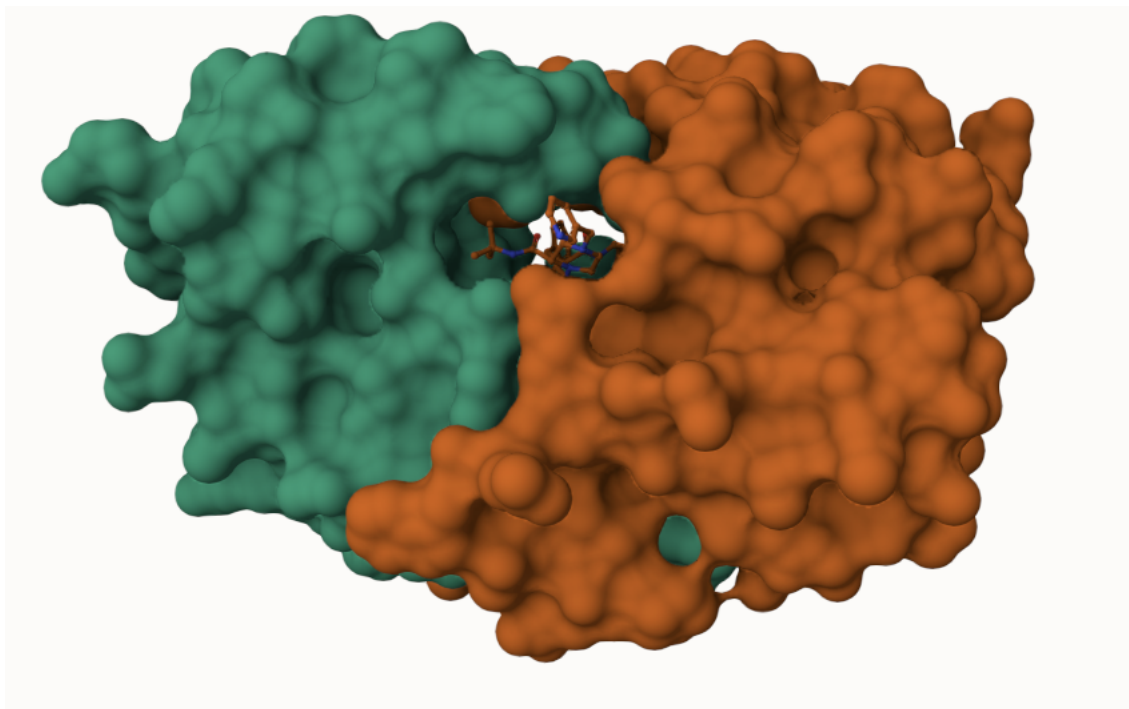
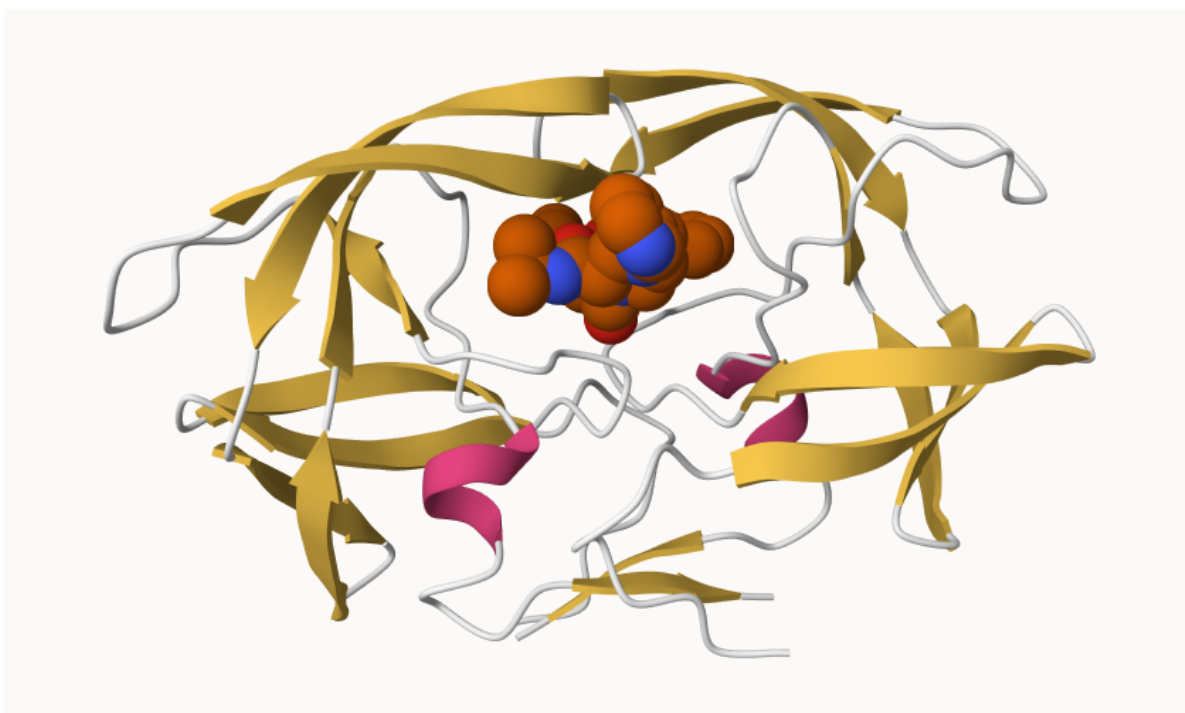


Figure 1: My first image of HIV-Pr with surface display showing ligand binding

A new clean image showing the catalytic ASP25 amino acids in both chains of the HIV-PR dimer along with the inhibitor and all important active site water.



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 because each of the molecules is being represented just by the oxygen atom.

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 critical water molecule would be H₂O with a residue number of 301.

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?



They could enter the binding site through the transient opening of the flaps that are covering the active site.

Bio3D package for structural bioinformatics

```
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?

H2O (water)

Q9: How many protein chains are in this structure?

2 protein chains, A & B

```
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

- would not ever download :(

Predicting functional motions of a single structure

Read an ADK structure from the PDB database:

```
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:

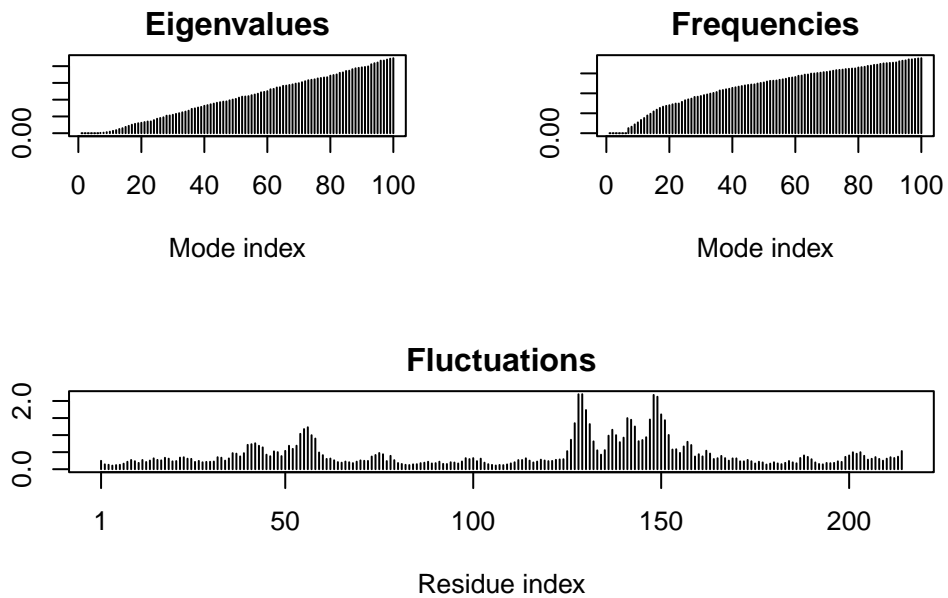
```
MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMLRAAVKSGSELGKQAKDIMDAGKLV
DELVIALVKERIAQEDCRNGFLLDGFPRPTIPQADAMKEAGINVDYVLEFDVPELIVDKI
VGRRVHAPSGRVYHVKFNPPKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQMTAPLIG
YYSKEAEAGNTKYAKVDGTPVAEVRADLEKILG
```

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

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

```
Building Hessian...      Done in 0.012 seconds.  
Diagonalizing Hessian... Done in 0.056 seconds.
```

```
plot(m)
```



write out our results as a wee trajectory/movie of predicted motions:

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




```
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 long

Align and superpose structures

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

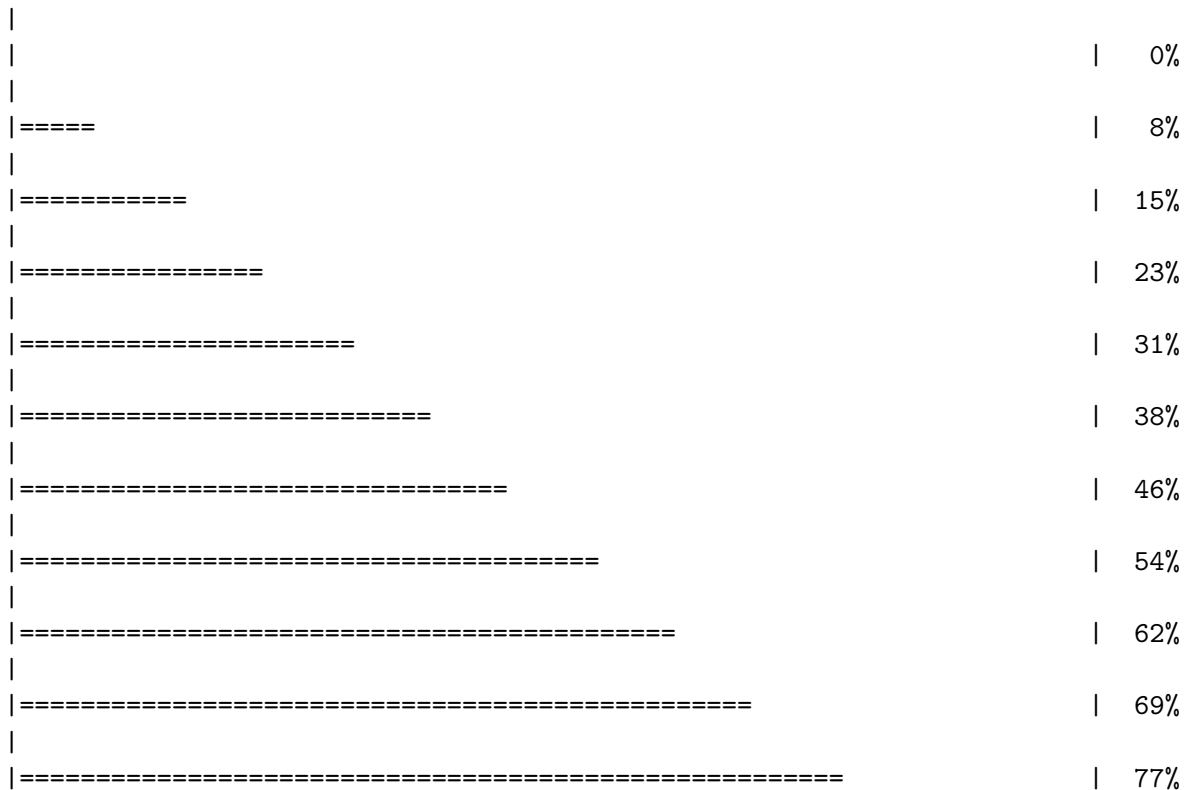
Warning in get.pdb(hits\$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/6HAP.pdb.gz exists. Skipping download

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



```

|
|=====| 85%
|
|=====| 92%
|
|=====| 100%

```

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

Reading PDB files:

```

pdbbs/split_chain/1AKE_A.pdb
pdbbs/split_chain/6S36_A.pdb
pdbbs/split_chain/6RZE_A.pdb
pdbbs/split_chain/3HPR_A.pdb
pdbbs/split_chain/1E4V_A.pdb
pdbbs/split_chain/5EJE_A.pdb
pdbbs/split_chain/1E4Y_A.pdb
pdbbs/split_chain/3X2S_A.pdb
pdbbs/split_chain/6HAP_A.pdb
pdbbs/split_chain/6HAM_A.pdb
pdbbs/split_chain/4K46_A.pdb
pdbbs/split_chain/3GMT_A.pdb
pdbbs/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: pdbbs/split_chain/1AKE_A.pdb
  PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 2  name: pdbbs/split_chain/6S36_A.pdb
  PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 3  name: pdbbs/split_chain/6RZE_A.pdb
  PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 4  name: pdbbs/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

```

```
ids <- basename.pdb(pdb$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"

```

```
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	CL (3),NA,MG (2)
6RZE_A	NA (3),CL (2)
3HPR_A	AP5
1E4V_A	AP5
5EJE_A	AP5,CO
1E4Y_A	AP5
3X2S_A	JPY (2),AP5,MG
6HAP_A	AP5
6HAM_A	AP5
4K46_A	ADP,AMP,PO4
3GMT_A	SO4 (2)
4PZL_A	CA,FMT,GOL

ligandName

1AKE_A BIS(ADENOSINE)-5'-
 PENTAPHOSPHATE
 6S36_A CHLORIDE ION (3),SODIUM ION,MAGNESIUM ION (2)
 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 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
 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 0139:H28 str. E24377A
 1E4Y_A Escherichia coli
 3X2S_A Escherichia coli str. K-12 substr. MDS42
 6HAP_A Escherichia coli 0139: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 INHIB
 6S36_A
 6RZE_A
 3HPR_A
 1E4V_A
 loop
 5EJE_A
 1E4Y_A
 loop

Cryst

3X2S_A
 conjugated adenylate kinase
 6HAP_A
 6HAM_A
 4K46_A
 3GMT_A
 4PZL_A

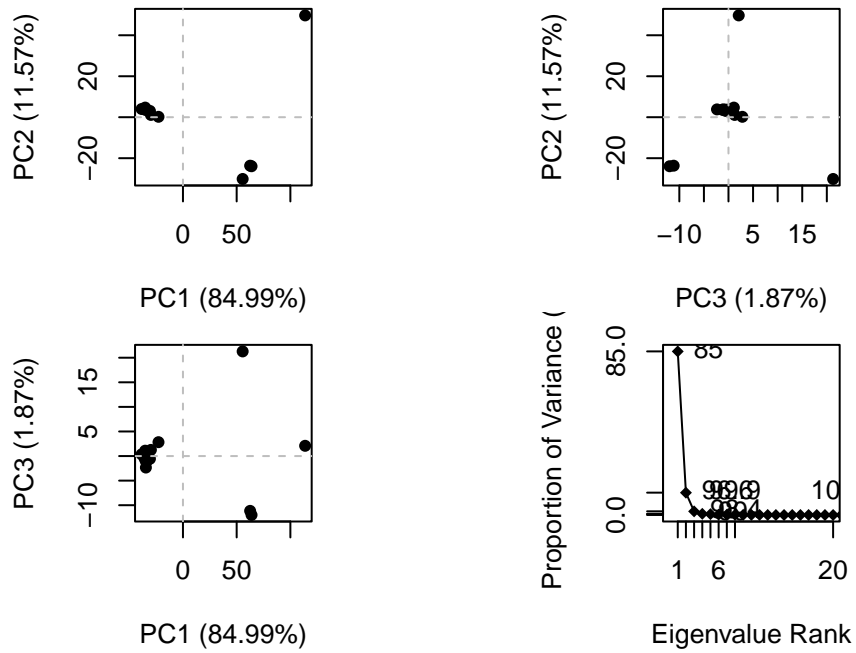
The crys

	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

Principal component analysis

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



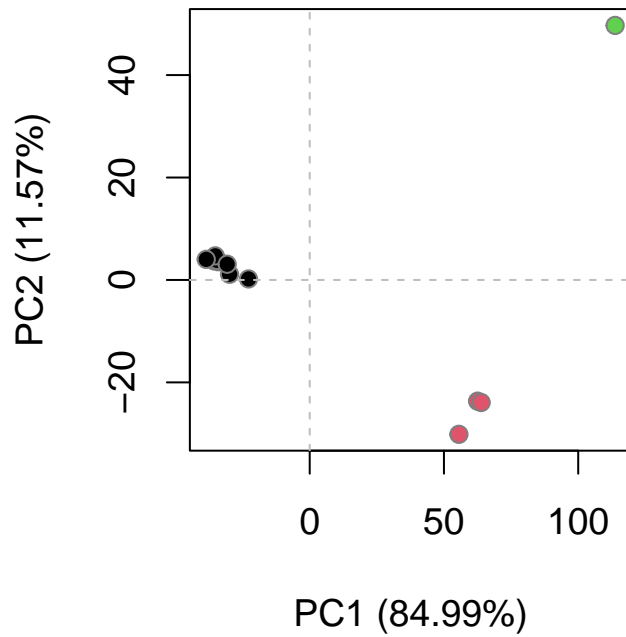
```
rd <- rmsd(pdfs)
```

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

```
hc.rd <- hclust(dist(rd))
```

```
grps.rd <- cutree(hc.rd, k=3)
```

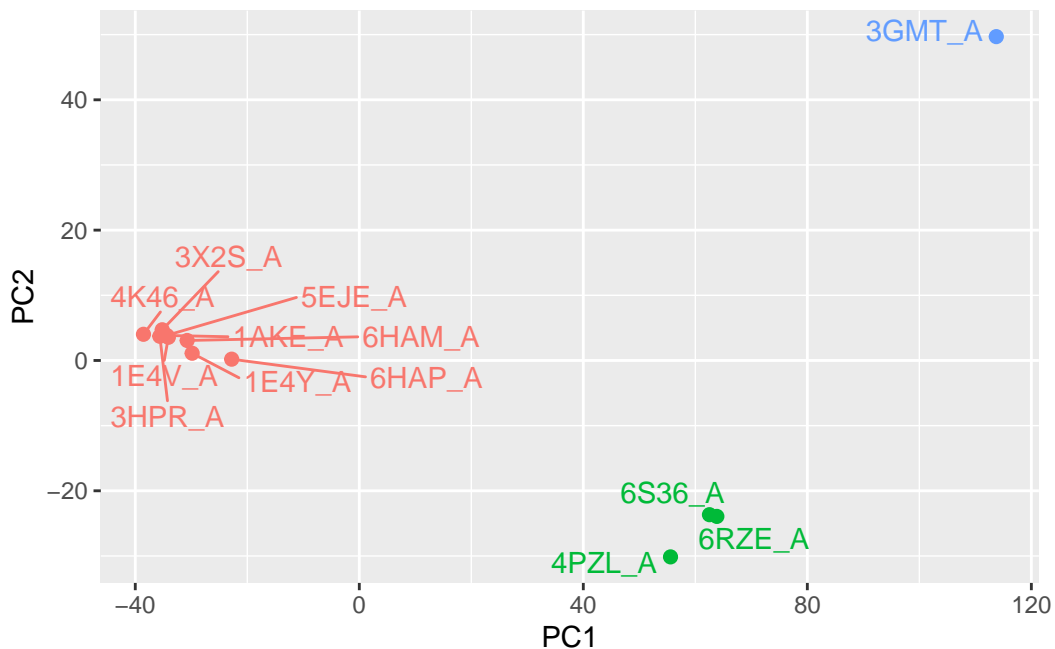
```
plot(pc.xray, 1:2, col="grey50", bg=grps.rd, pch=21, cex=1)
```



```
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
```



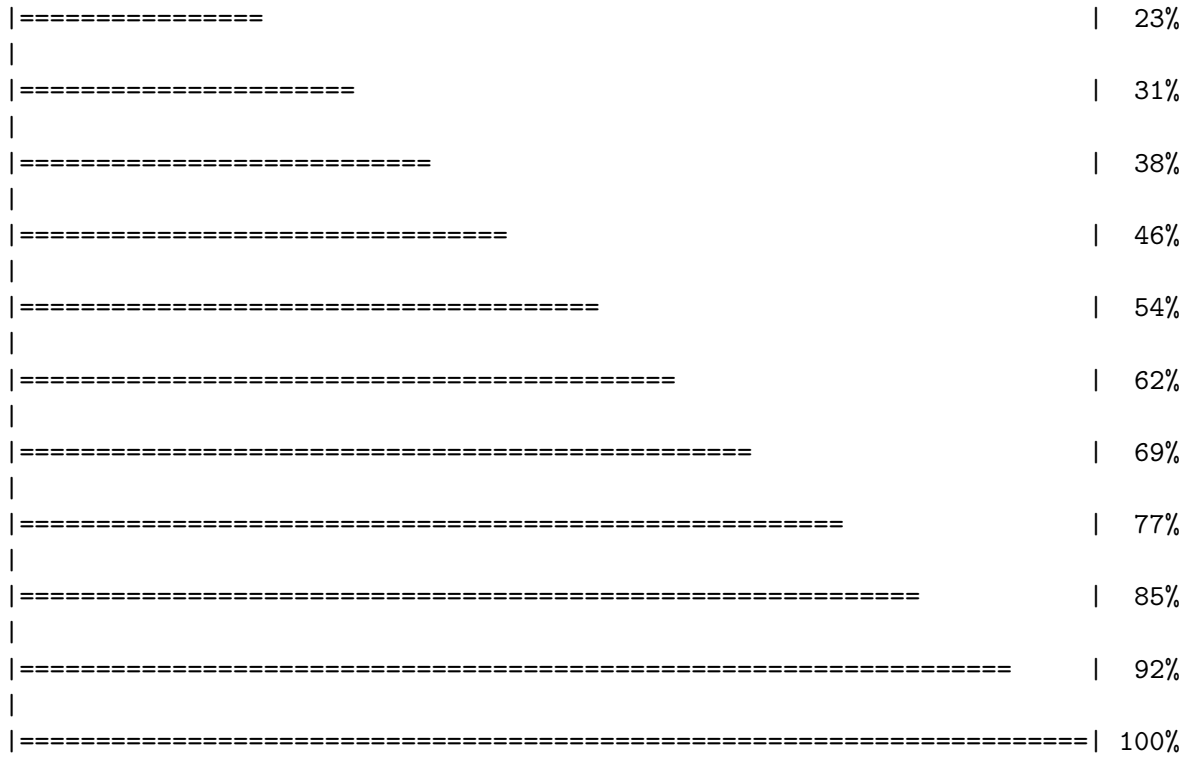
Normal mode analysis [optional]

```
modes <- nma(pdb)
```

Details of Scheduled Calculation:

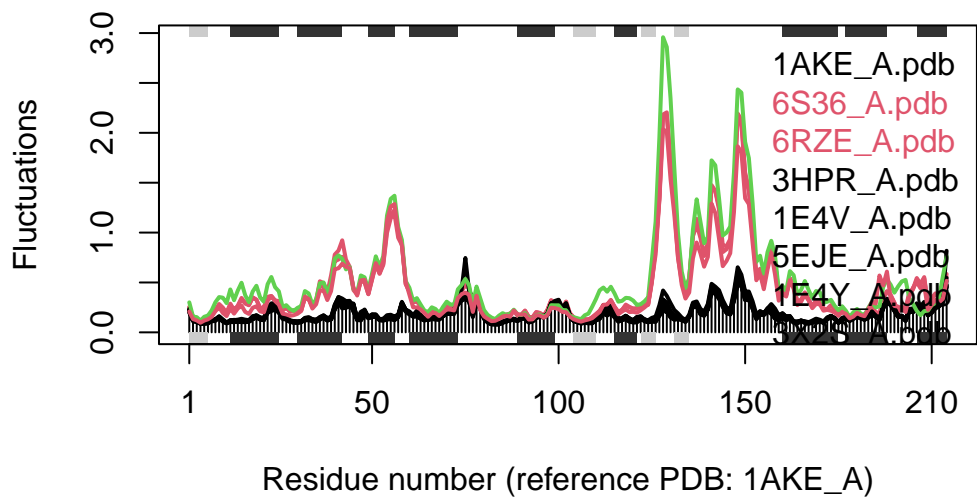
```
... 13 input structures
... storing 606 eigenvectors for each structure
... dimension of x$U.subspace: ( 612x606x13 )
... coordinate superposition prior to NM calculation
... aligned eigenvectors (gap containing positions removed)
... estimated memory usage of final 'eNMA' object: 36.9 Mb
```

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



```
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?

I notice that the black and colored lines have the same overall pattern but they differ in amplitude. The biggest differences can be seen in the binding regions.