

Class 10

AUTHOR
Erin

Q1: What percentage of structures in the PDB are solved by X-Ray and Electron Microscopy. 93.15962 84.83% is X.ray while 8.33% is EM.

```
PDB.df <- read.csv("Data Export Summary.csv", row.names=1)
PDB.df$X.ray <- as.numeric(gsub(",", "", PDB.df$X.ray))
PDB.df$EM <- as.numeric(gsub(",", "", PDB.df$EM))
PDB.df$NMR <- as.numeric(gsub(",", "", PDB.df$NMR))
head(PDB.df)
```

	X.ray	EM	NMR	Multiple.methods	Neutron	Other
Protein (only)	158844	11759	12296	197	73	32
Protein/Oligosaccharide	9260	2054	34	8	1	0
Protein/NA	8307	3667	284	7	0	0
Nucleic acid (only)	2730	113	1467	13	3	1
Other	164	9	32	0	0	0
Oligosaccharide (only)	11	0	6	1	0	4
Total						
Protein (only)	183,201					
Protein/Oligosaccharide	11,357					
Protein/NA	12,265					
Nucleic acid (only)	4,327					
Other	205					
Oligosaccharide (only)	22					

```
total_Xray_EM <- sum(PDB.df$X.ray) + sum(PDB.df$EM)
grand_total <- sum(PDB.df$X.ray, PDB.df$EM, PDB.df$NMR, PDB.df$Multiple.methods, PDB.df$Neutron, PDB.df$Other)

percentage_Xray_EM <- (total_Xray_EM / grand_total) * 100
percentage_Xray_EM
```

[1] 93.15962

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

	X.ray	EM	NMR	Multiple.methods	Neutron	Other
Protein (only)	158,844	11,759	12,296	197	73	32
Protein/Oligosaccharide	9,260	2,054	34	8	1	0
Protein/NA	8,307	3,667	284	7	0	0
Nucleic acid (only)	2,730	113	1,467	13	3	1
Other	164	9	32	0	0	0
Oligosaccharide (only)	11	0	6	1	0	4
Total						
Protein (only)	183,201					
Protein/Oligosaccharide	11,357					
Protein/NA	12,265					
Nucleic acid (only)	4,327					
Other	205					
Oligosaccharide (only)	22					

```
# create working snippet
x <- stats$X.ray
x
```

[1] "158,844" "9,260" "8,307" "2,730" "164" "11"

```
as.numeric(gsub(",", "", x))
```

```
[1] 158844 9260 8307 2730 164 11
```

```
rm.comma <- function(x){
  as.numeric(gsub(",", "", x))
}
rm.comma(stats$X.ray)
```

```
[1] 158844 9260 8307 2730 164 11
```

```
pdbstats <- apply(stats, 2, rm.comma)
rownames(pdbstats) <- rownames(stats)
head(pdbstats)
```

	X.ray	EM	NMR	Multiple.methods	Neutron	Other
Protein (only)	158844	11759	12296	197	73	32
Protein/Oligosaccharide	9260	2054	34	8	1	0
Protein/NA	8307	3667	284	7	0	0
Nucleic acid (only)	2730	113	1467	13	3	1
Other	164	9	32	0	0	0
Oligosaccharide (only)	11	0	6	1	0	4
Total						
Protein (only)	183201					
Protein/Oligosaccharide	11357					
Protein/NA	12265					
Nucleic acid (only)	4327					
Other	205					
Oligosaccharide (only)	22					

```
pdbtotals <- apply(pdbstats, 2, sum)
pdbtotals
```

X.ray	EM	NMR	Multiple.methods
179316	17602	14119	226
Neutron	Other	Total	
77	37	211377	

```
## solved by different methods
round(pdbtotals / pdbtotals["Total"]*100, 2)
```

X.ray	EM	NMR	Multiple.methods
84.83	8.33	6.68	0.11
Neutron	Other	Total	
0.04	0.02	100.00	

Q2: What proportion of structures in the PDB are protein? 86.67%

```
# Step 1: Use rowSums() to get the total count for each Molecular Type
PDB.df$total <- rowSums(PDB.df[, c('X.ray', 'EM', 'NMR', 'Multiple.methods', 'Neutron', 'Other')], na.rm = TRUE)

# Step 2: Extract the total count for 'Protein (only)' which is the first row in the dataset
protein_total <- PDB.df$total[1]

# Step 3: Sum the total count for all molecular types to get the grand total
grand_total <- sum(PDB.df$total)

# Step 4: Calculate the proportion of protein structures
protein_proportion <- protein_total / grand_total
protein_proportion
```

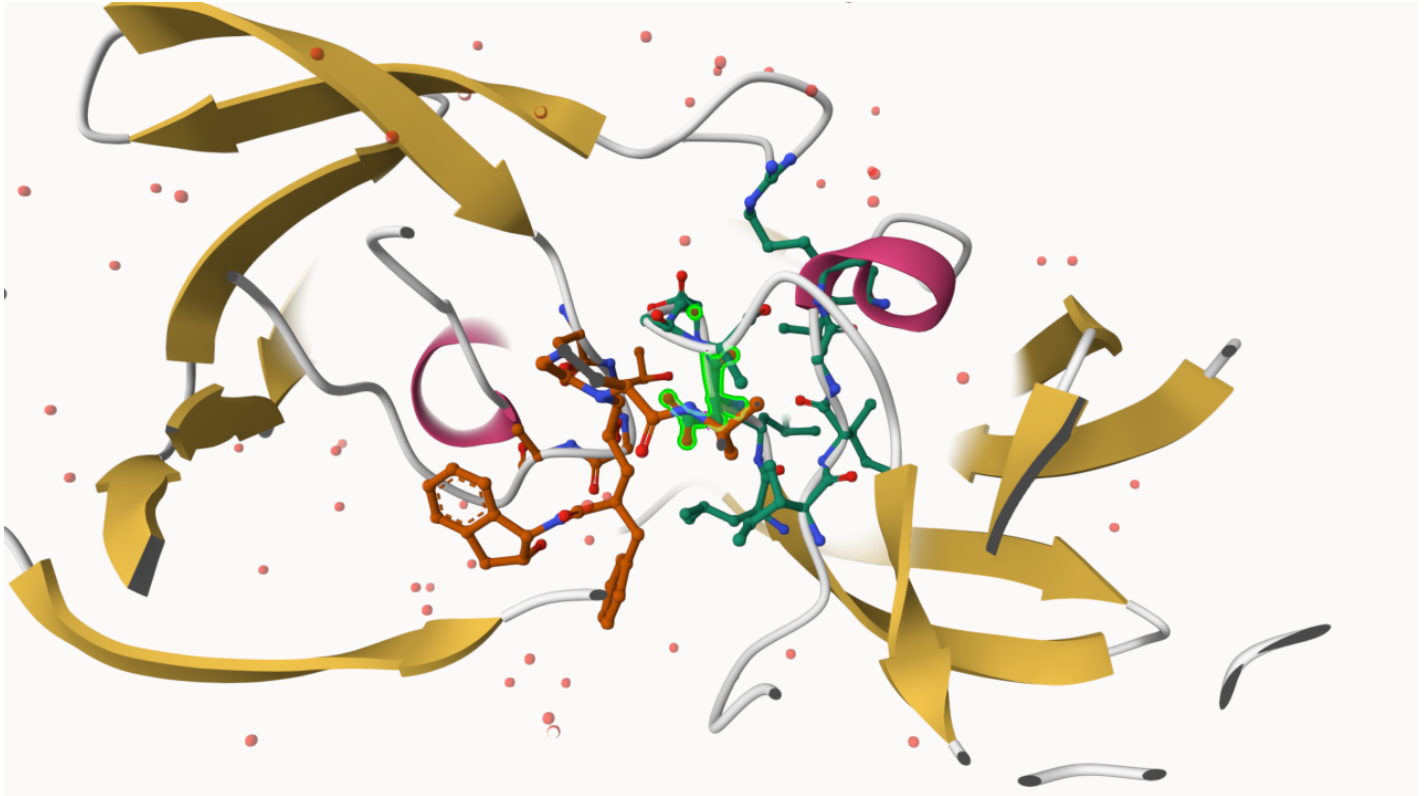
[1] 0.8667026

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? 7434

Q4: Water molecules normally have 3 atoms. Why do we see just one atom per water molecule in this structure? The resolution limit is only 2.00Å, therefore hydrogen is too small to be resolved.

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 Water molecule is identifiable. Water 308 is responsible for stabilizing the ligand-protein interaction by H bond.

Q6



Q7: [Optional] 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 can you identify secondary structure elements that are likely to only form in the dimer rather than the monomer?

```
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
VNIIGRNLLTQIGCTLNFP
```

+ 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, MK1 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	elenu	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	elenu	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>

```
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
DELVIALVKERIAQEDCRNGFLDGFPRTPQADAMKEAGINVDYVLEFDVPDELIVDKI
VGRRVHAPSGRVYHVKNPPKVEGKDDVTGEELTRKDDQEETVRKRLVEYHQMTAPLIG
```

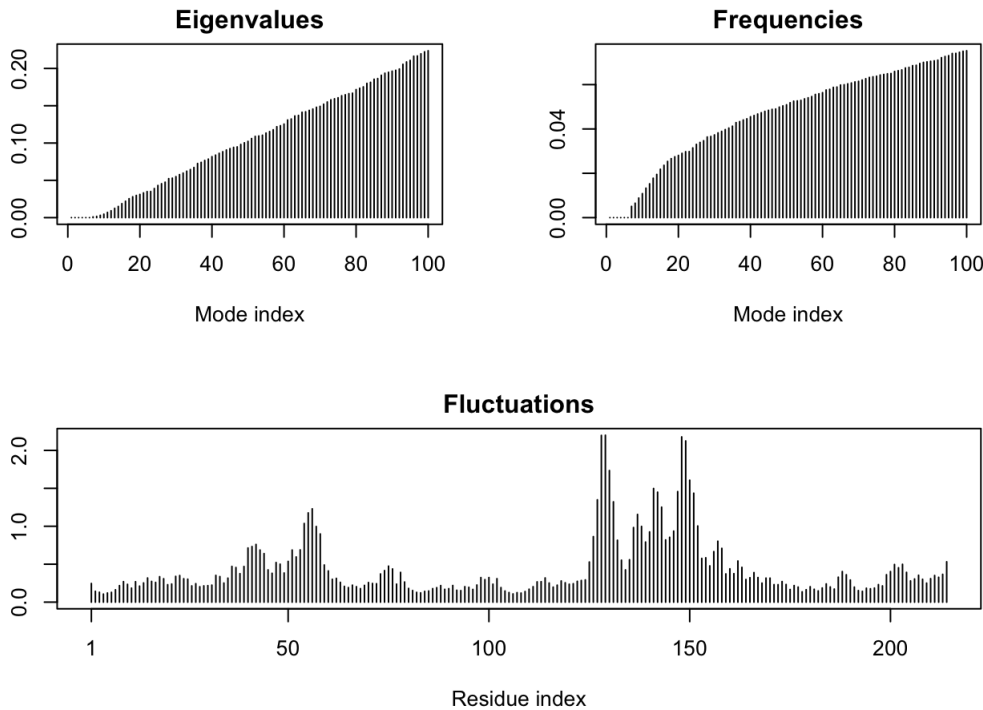
YYSKEAEAGNTKYAKVDGTPVAEVRADLEKILG

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

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

```
Building Hessian...      Done in 0.075 seconds.  
Diagonalizing Hessian... Done in 0.26 seconds.
```

```
plot(m)
```



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

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?: Grantlab/bio3d-view Q12. True or False? Functions from the devtools package can be used to install packages from GitHub and BitBucket? True

```
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 MRILLGAPGAGKGTQAQFIMEKYGIPQISTGDMRLRAAVKSGSELGKQAKDIMDAGKLV  
      1      .      .      .      .      .      60  
  
     61      .      .      .      .      .      120  
pdb|1AKE|A DELVIALVKERIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFDVPDELIVDRI
```

```

        61      .      .      .      .      .      .      120
        121     .      .      .      .      .      .      180
pdb|1AKE|A  VGRRVHAPSGRVYHVKNPPKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQMTAPLIG
        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 = N4VFDRR6016

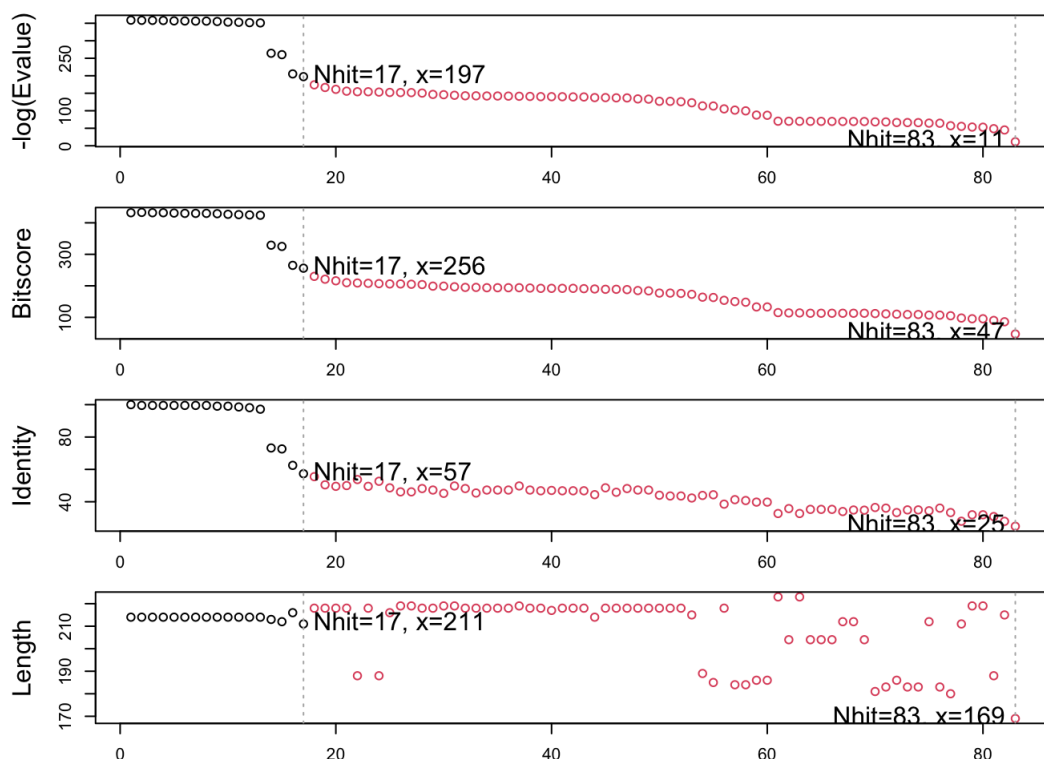
.

Reporting 83 hits

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

```
* Possible cutoff values: 197 11
    Yielding Nhits:      17 83
```

```
* Chosen cutoff value of: 197
    Yielding Nhits:      17
```



```
# List out some 'top hits'
head(hits$pdb.id)
```

```
[1] "1AKE_A" "8BQF_A" "4X8M_A" "6S36_A" "6RZE_A" "4X8H_A"
```

```
hits <- NULL
```

```
hits$pdb.id <- c('1AKE_A', '6S36_A', '6RZE_A', '3HPR_A', '1E4V_A', '5EJE_A', '1E4Y_A', '3X2S_A', '6HAP_A', '6HAM_A', '4K46_
```

```
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
```

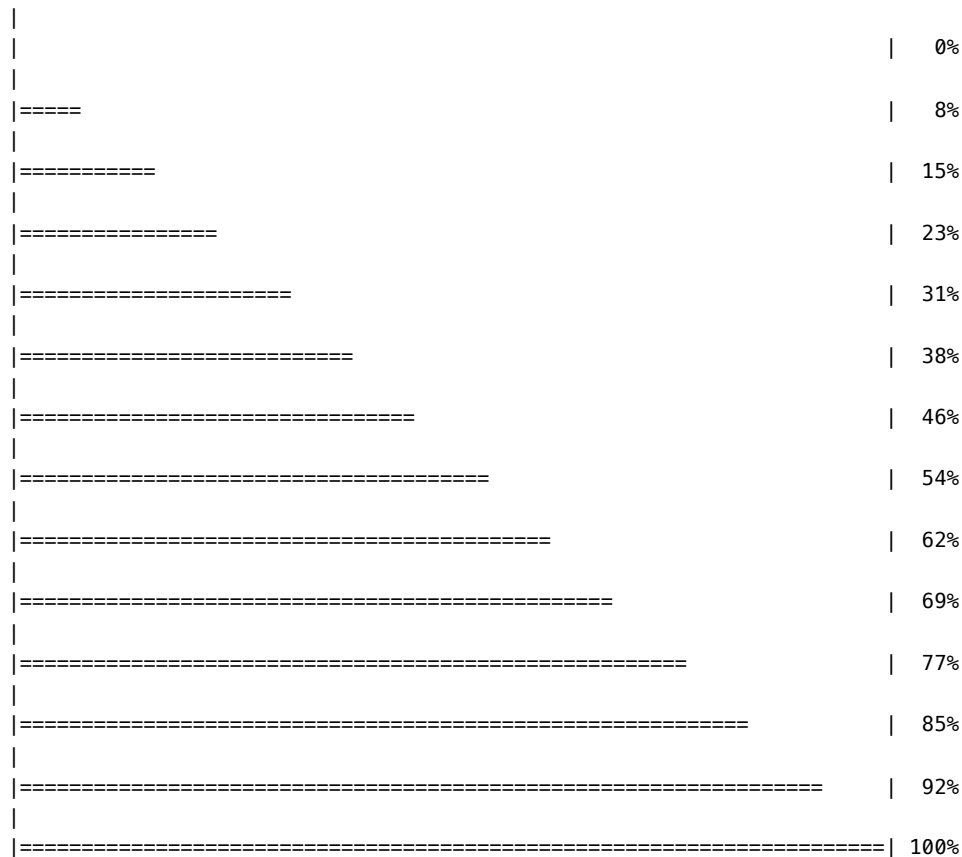
```
Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/6HAP.pdb.gz exists. Skipping download
```

```
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
```



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

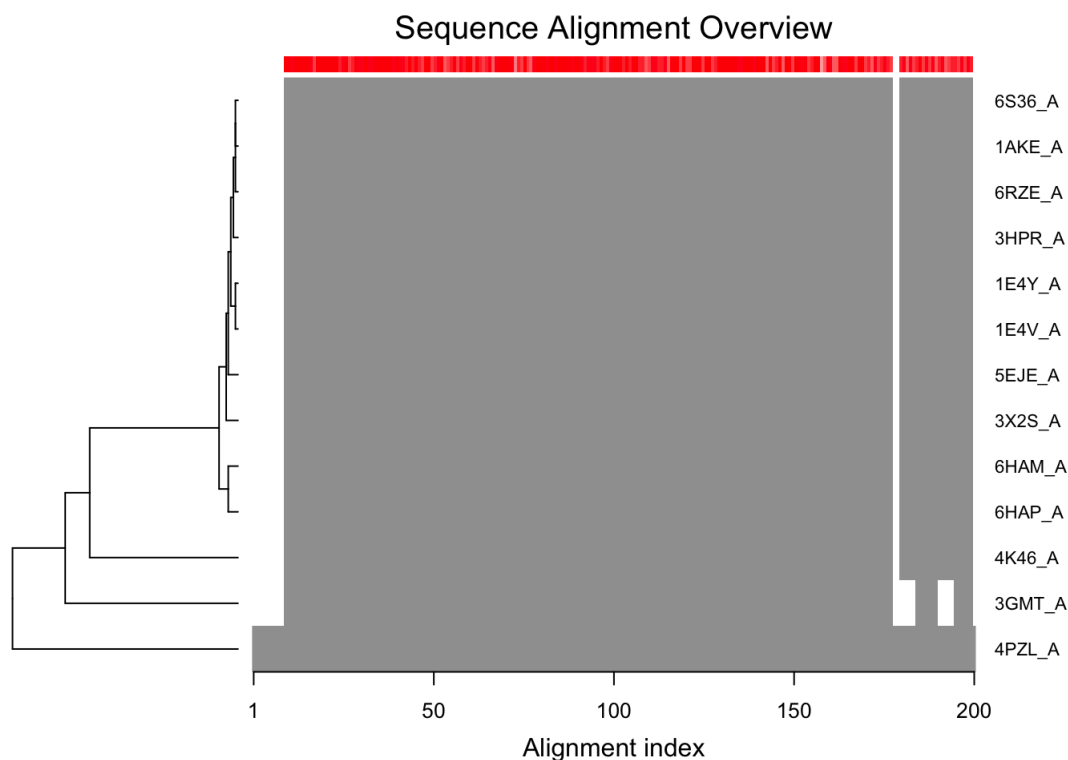
```

```

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

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

```



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

```
[1] "Escherichia coli"
[2] "Escherichia coli K-12"
[3] "Escherichia coli 0139: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, active site lid (ADK_lid)		
	6RZE_A	1.69	<NA>	Adenylate kinase, active site lid (ADK_lid)		
	3HPR_A	2.00	<NA>	Adenylate kinase, active site lid (ADK_lid)		
	1E4V_A	1.85	Adenylate kinase	Adenylate kinase, active site lid (ADK_lid)		
	5EJE_A	1.90	<NA>	Adenylate kinase, active site lid (ADK_lid)		
	1E4Y_A	1.85	Adenylate kinase	Adenylate kinase, active site lid (ADK_lid)		
	3X2S_A	2.80	<NA>	Adenylate kinase, active site lid (ADK_lid)		
	6HAP_A	2.70	<NA>	Adenylate kinase, active site lid (ADK_lid)		
	6HAM_A	2.55	<NA>	Adenylate kinase, active site lid (ADK_lid)		
	4K46_A	2.01	<NA>	Adenylate kinase, active site lid (ADK_lid)		
	3GMT_A	2.10	<NA>	Adenylate kinase, active site lid (ADK_lid)		
	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,C0				
	1E4Y_A	AP5				
	3X2S_A	JPY (2),AP5,MG				
	6HAP_A	AP5				
	6HAM_A	AP5				
	4K46_A	ADP,AMP,P04				
	3GMT_A	S04 (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

structureTitle

1AKE_A STRUCTURE OF THE COMPLEX BETWEEN ADENYLATE KINASE FROM ESCHERICHIA COLI AND THE INHIBITOR AP5A REFINED AT 1.9 ANGSTROMS RESOLUTION: A MODEL FOR A CATALYTIC TRANSITION STATE

6S36_A

Crystal structure of E. coli Adenylate kinase R119K mutant

6RZE_A

Crystal structure of E. coli Adenylate kinase R119A mutant

3HPR_A

structure of V148G adenylate kinase from E. coli, in complex with Ap5A

1E4V_A

Mutant G10V of adenylate kinase from E. coli, modified in the Gly-loop

5EJE_A

coli Adenylate kinase G56C/T163C double mutant in complex with Ap5a

1E4Y_A

Mutant P9L of adenylate kinase from E. coli, modified in the Gly-loop

3X2S_A

Crystal structure of pyrene-conjugated adenylate kinase

6HAP_A

Adenylate kinase

6HAM_A

Adenylate kinase

4K46_A

Crystal Structure of Adenylate Kinase from Photobacterium profundum

3GMT_A

Crystal structure of adenylate kinase from burkholderia pseudomallei

4PZL_A

adenylate kinase from Francisella tularensis subsp. tularensis SCHU S4

Crystal

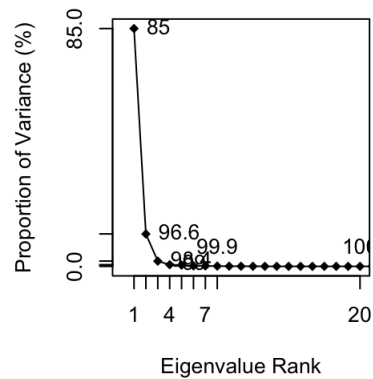
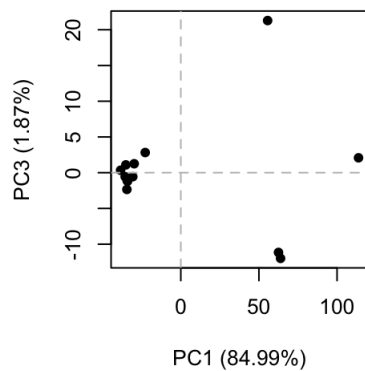
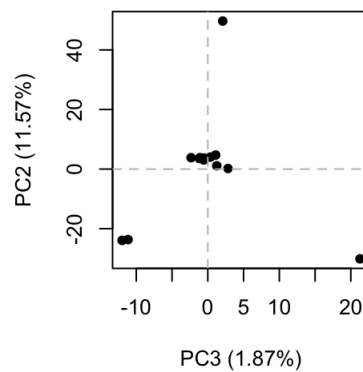
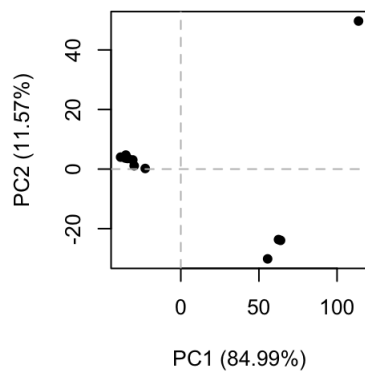
Crystal structure of E.

The crystal structure of

		citation	rObserved	rFree
1AKE_A	Muller, C.W., et al. J Mol Biol (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	space	Group
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

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

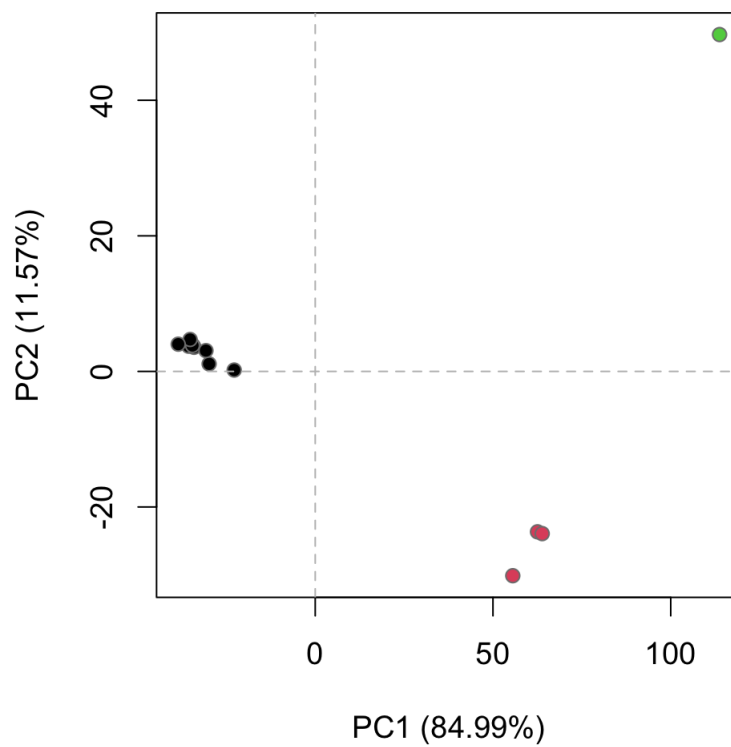


```
# Calculate RMSD
rd <- rmsd(pdbbs)
```

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)
```

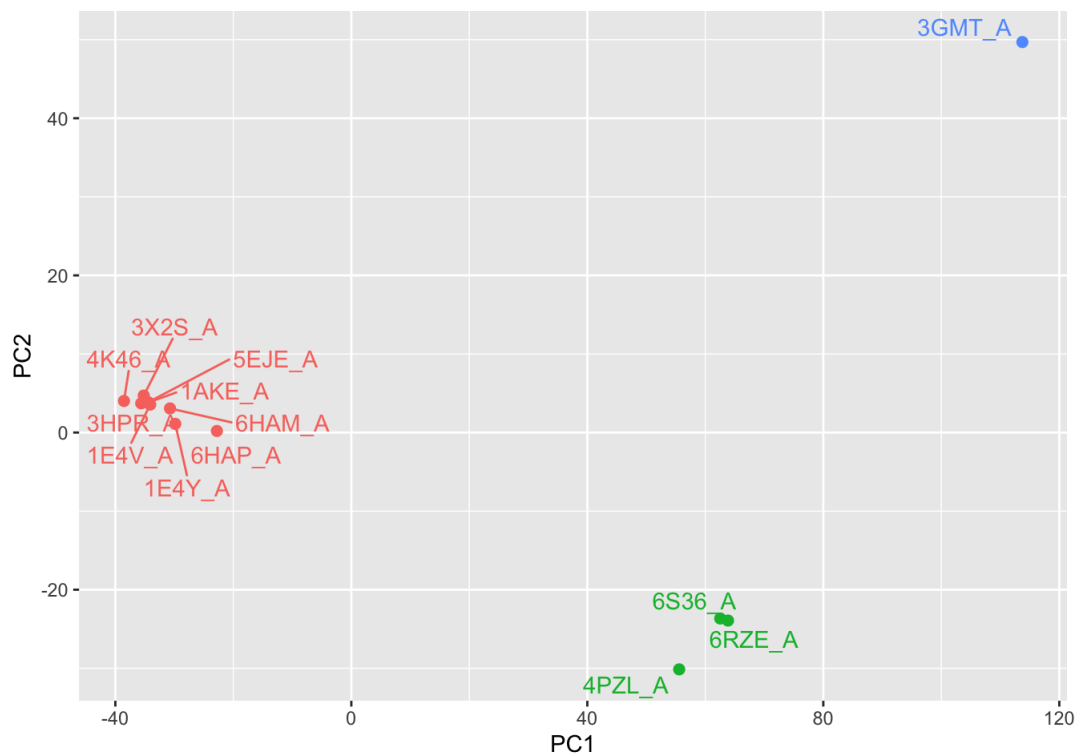


```
# Visualize first principal component
pc1 <- mkrj(pc.xray, pc=1, file="pc_1.pdb")
```

```
#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)

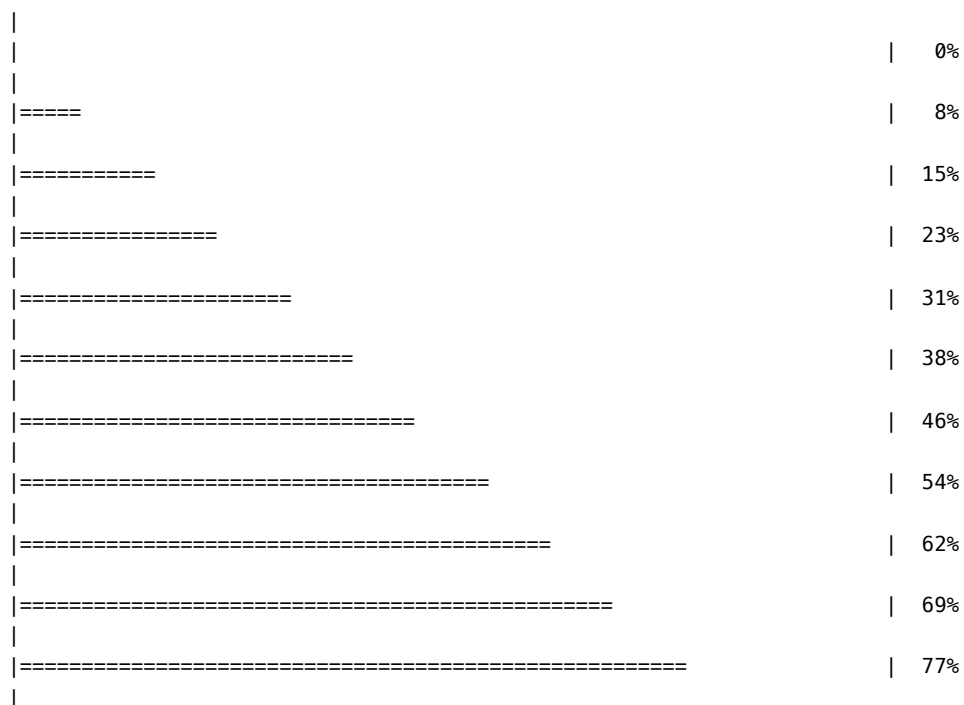
p <- ggplot(df) +
  aes(PC1, PC2, col=col, label=ids) +
  geom_point(size=2) +
  geom_text_repel(max.overlaps = 20) +
  theme(legend.position = "none")
p
```



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

Details of Scheduled Calculation:

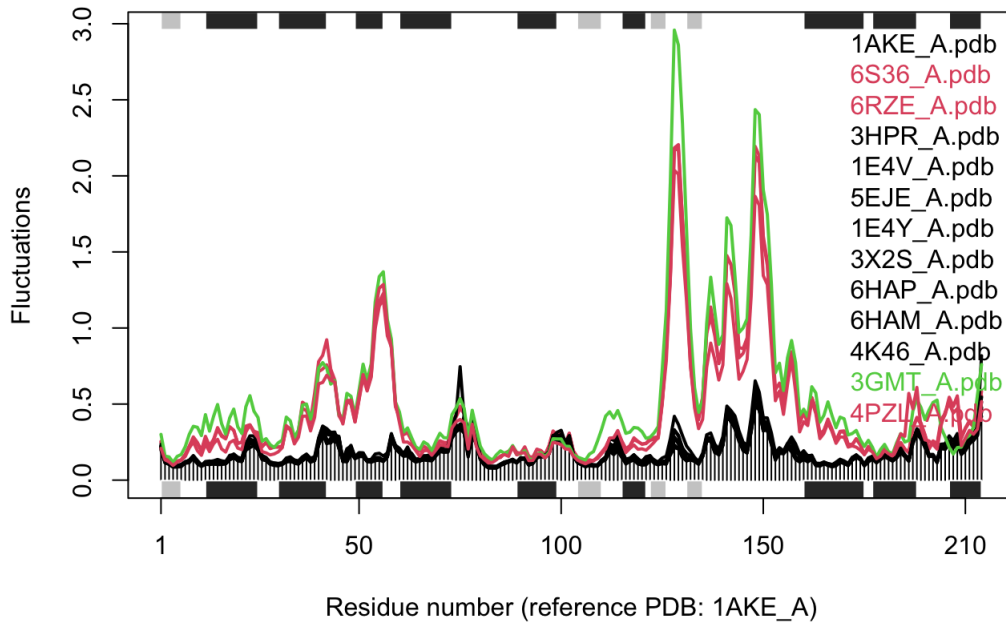
```
... 13 input structures
... storing 606 eigenvectors for each structure
... dimension of x$U.subspace: ( 612x606x13 )
... coordinate superposition prior to NM calculation
... aligned eigenvectors (gap containing positions removed)
... estimated memory usage of final 'eNMA' object: 36.9 Mb
```



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

```
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 colored lines represent different structures of the same protein (indicated by PDB codes like 1AKE_A, 4X8M_A, etc.), while the black line may represent the average or a reference structure. The black bars at the top indicate where the differences between the colored lines and the black line are statistically significant. Overall, the fluctuation patterns of the colored lines follow the same general trend as the black line. This suggests that the regions of flexibility and rigidity are relatively conserved across the different structures. The most significant differences appear to be at specific points where the colored lines show peaks that are much higher than the black line. These are likely regions where certain structures have more flexibility or exhibit more movement than the average or reference structure. This could be due to differences in crystal packing, the presence of bound ligands or other molecules, or mutations in some of the structures.