

class09

AUTHOR

Ziheng Huang

1: Introduction to the RCSB Protein Data Bank (PDB)

```
PDBSummary <- "PDBSummary.csv"
PDBSummary <- read.csv("PDBSummary.csv", row.names=1)
PDBSummary
```

| | X.ray | NMR | EM | Multiple.methods | Neutron | Other |
|-------------------------|---------|--------|-------|------------------|---------|-------|
| Protein (only) | 150,342 | 12,053 | 8,534 | 188 | 72 | 32 |
| Protein/Oligosaccharide | 8,866 | 32 | 1,540 | 6 | 0 | 0 |
| Protein/NA | 7,911 | 278 | 2,681 | 6 | 0 | 0 |
| Nucleic acid (only) | 2,510 | 1,425 | 74 | 13 | 2 | 1 |
| Other | 154 | 31 | 6 | 0 | 0 | 0 |
| Oligosaccharide (only) | 11 | 6 | 0 | 1 | 0 | 4 |
| Total | | | | | | |
| Protein (only) | 171,221 | | | | | |
| Protein/Oligosaccharide | 10,444 | | | | | |
| Protein/NA | 10,876 | | | | | |
| Nucleic acid (only) | 4,025 | | | | | |
| Other | 191 | | | | | |
| Oligosaccharide (only) | 22 | | | | | |

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

X-ray:

86.28665

EM:

6.522546

X-ray + EM:

92.80919

```
#total
tot <- sum(as.numeric(sub(',', '', PDBSummary$Total)))
# X-ray
XR = 100*sum(as.numeric(sub(',', '', PDBSummary$X.ray)))/tot
XR
```

[1] 86.28665

```
# EM:  
EM = 100*sum(as.numeric(sub('','',PDBSummary$EM)))/tot  
EM
```

[1] 6.522546

```
# X-ray + EM  
XR + EM
```

[1] 92.80919

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

0.9784631

```
Prot = (as.numeric(sub('','',PDBSummary["Protein (only)",]$Total)) + as.numeric(sub('','',  
Prot
```

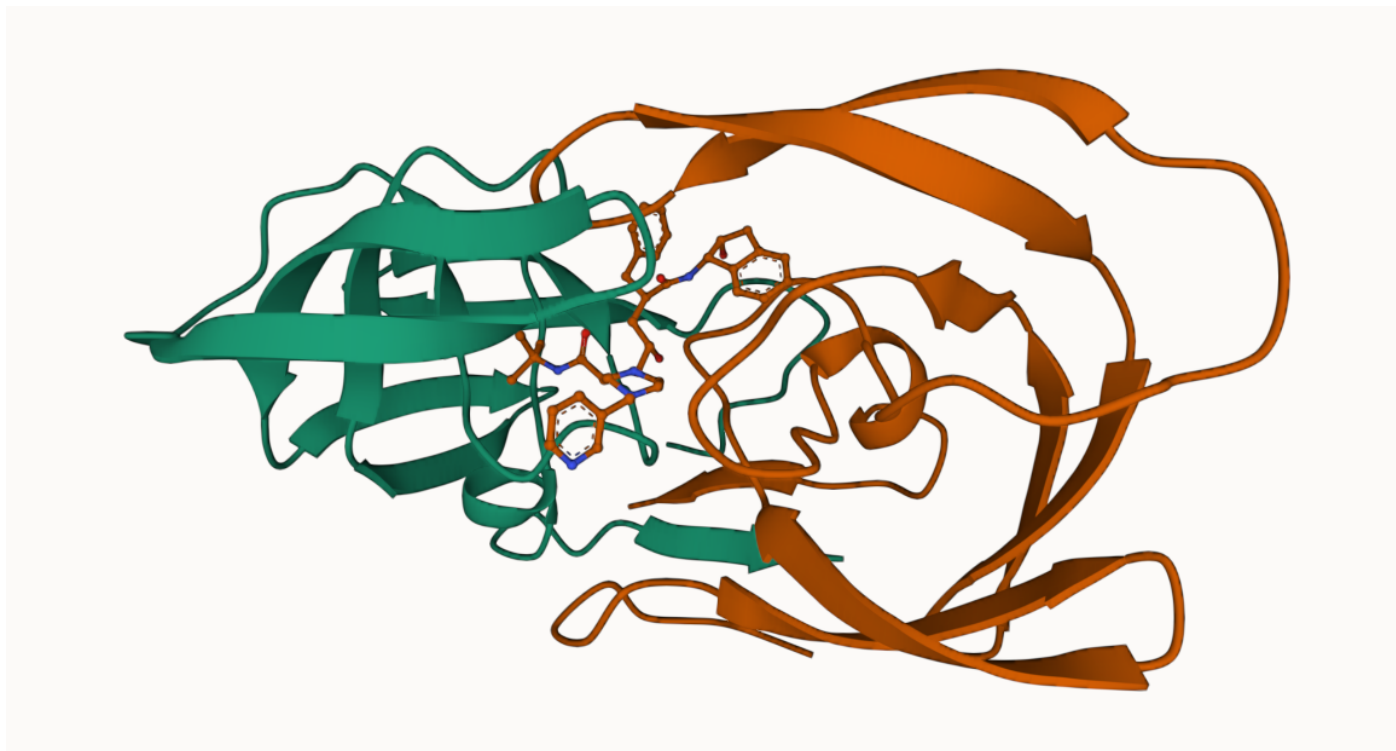
[1] 0.9784631

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?

Query "HIV" matches 4707 structures

2. Visualizing the HIV-1 protease structure

Using Mol*



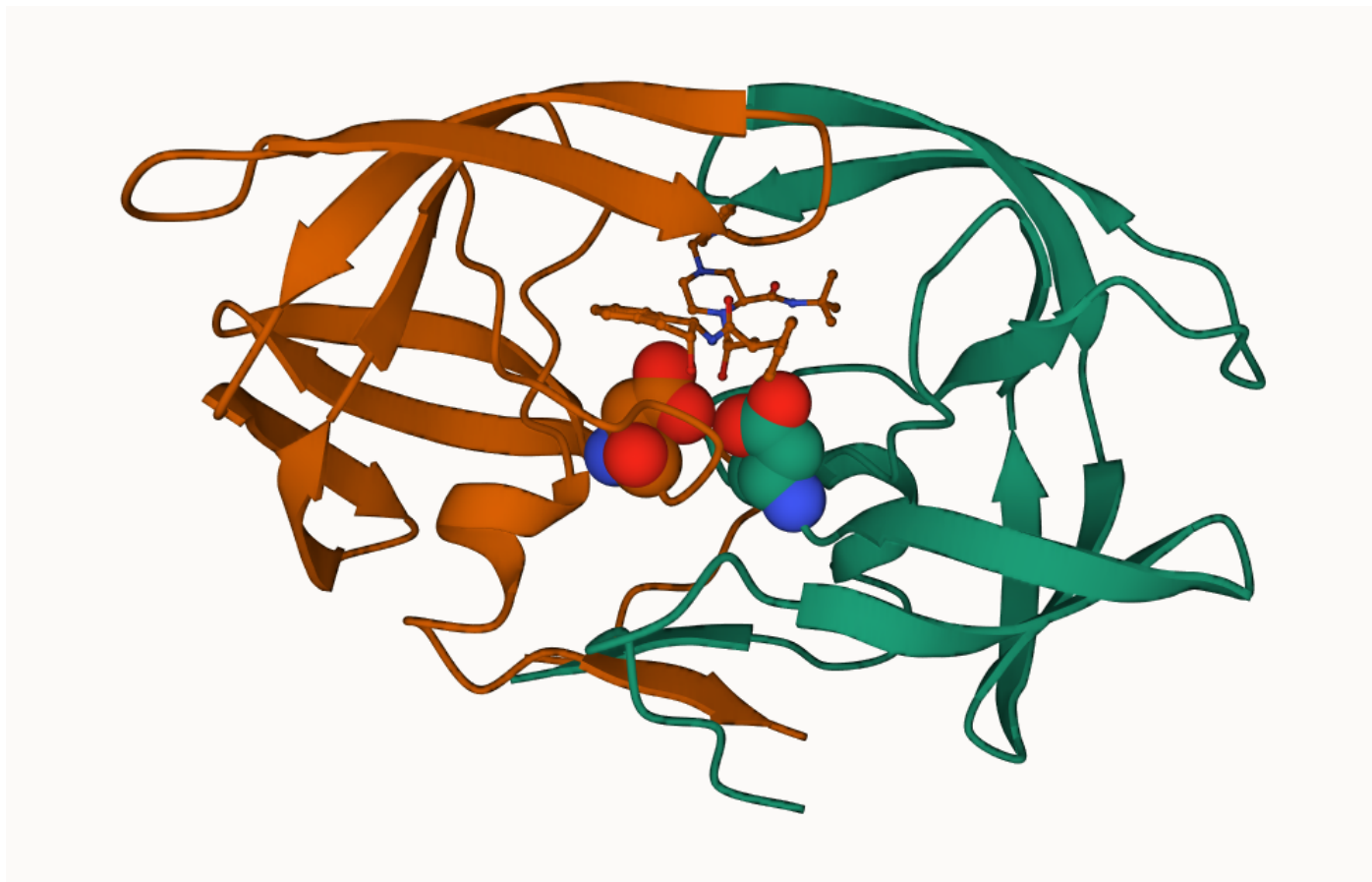
Q4: Water molecules normally have 3 atoms. Why do we see just one atom per water molecule in this structure?

hydrogen are too small, so we only see 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 ?

H2O 308

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 (we recommend "Ball & Stick" for these side-chains). Add this figure to your Quarto document.



3. Introduction to Bio3D in R

load data

```
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

Q8: Name one of the two non-protein residues? MK1

Q9: How many protein chains are in this structure? 2

Atom records of a PDB file:

```
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

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("lake_A"): Removing existing file: seqs.fasta

Fetching... Please wait. Done.

aa

```

      1      .      .      .      .      .      .      60
pdb|1AKE|A  MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMRLRAAVKSGSELGKQAKDIMDAGKLV
      1      .      .      .      .      .      .      60

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

     121      .      .      .      .      .      .      180
pdb|1AKE|A  VGRRVHAPSGRVYHVKNPPKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQM
     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)
```

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

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

Use these for analysis:

```
hits <- NULL
hits$ pdb.id <- c('1AKE A','6S36 A','6RZE A','3HPR A','1E4V A','5EJE A','1E4Y A','3X2S A',
```

Download files

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

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

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

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

Warning in get.pdb(hits\$ pdb.id, path = "pdbc", split = TRUE, gzip = TRUE): pdbc/
3HPR.pdb.gz exists. Skipping download

Warning in get.pdb(hits\$ pdb.id, path = "pdbc", split = TRUE, gzip = TRUE): pdbc/
1E4V.pdb.gz exists. Skipping download

Warning in get.pdb(hits\$ pdb.id, path = "pdbc", split = TRUE, gzip = TRUE): pdbc/
5EJE.pdb.gz exists. Skipping download

Warning in get.pdb(hits\$ pdb.id, path = "pdbc", split = TRUE, gzip = TRUE): pdbc/
1E4Y.pdb.gz exists. Skipping download

Warning in get.pdb(hits\$ pdb.id, path = "pdbc", split = TRUE, gzip = TRUE): pdbc/
3X2S.pdb.gz exists. Skipping download

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

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

Warning in get.pdb(hits\$ pdb.id, path = "pdbc", split = TRUE, gzip = TRUE): pdbc/
4K46.pdb.gz exists. Skipping download

Warning in get.pdb(hits\$ pdb.id, path = "pdbc", split = TRUE, gzip = TRUE): pdbc/
3GMT.pdb.gz exists. Skipping download

Warning in get.pdb(hits\$ pdb.id, path = "pdbc", split = TRUE, gzip = TRUE): pdbc/
4PZL.pdb.gz exists. Skipping download

| | | |
|-------|--|-----|
| | | |
| | | 0% |
| | | |
| ===== | | 8% |
| | | |
| ===== | | 15% |

| | | |
|-------|--|------|
| | | |
| ===== | | 23% |
| | | |
| ===== | | 31% |
| | | |
| ===== | | 38% |
| | | |
| ===== | | 46% |
| | | |
| ===== | | 54% |
| | | |
| ===== | | 62% |
| | | |
| ===== | | 69% |
| | | |
| ===== | | 77% |
| | | |
| ===== | | 85% |
| | | |
| ===== | | 92% |
| | | |
| ===== | | 100% |

Align and superpose structures

```
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
```

```
...
```

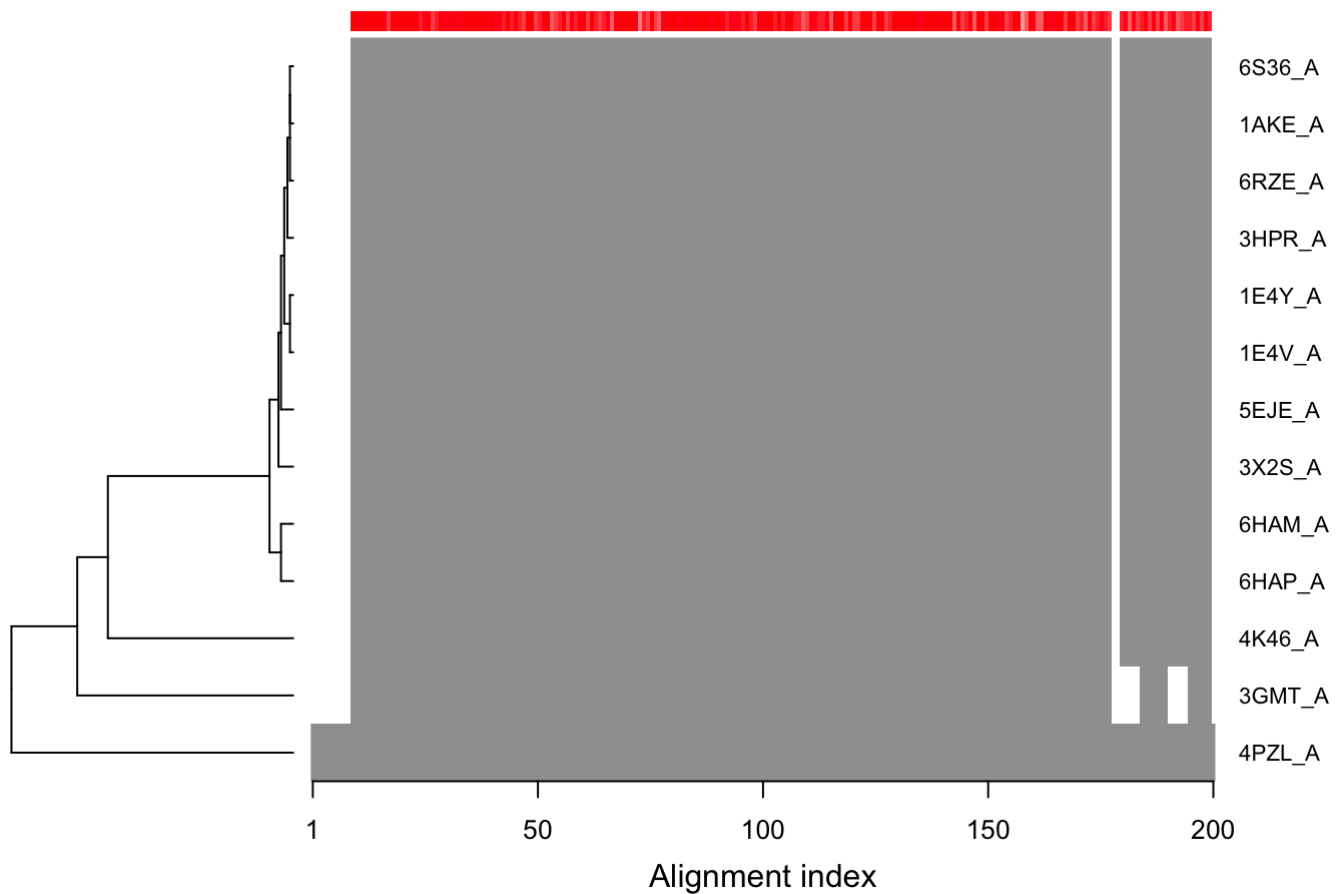

Extracting sequences

```
pdb/seq: 1  name: 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)

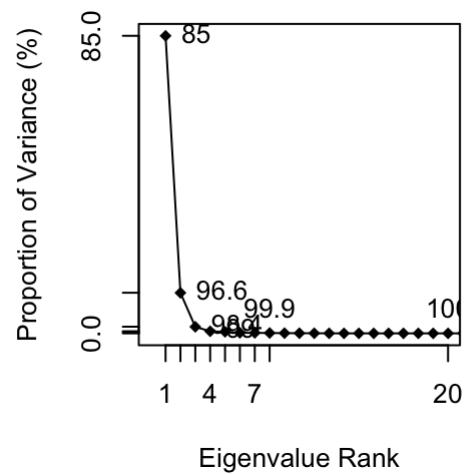
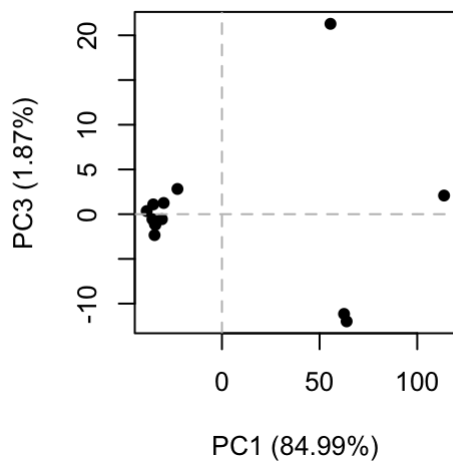
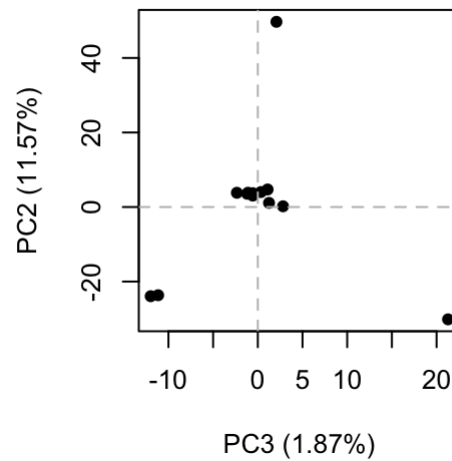
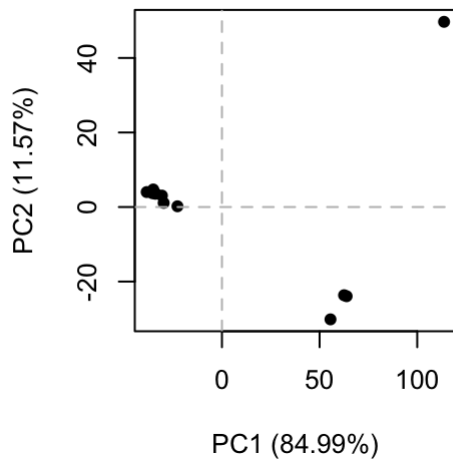
# Draw schematic alignment
plot(pdb, labels=ids)
```

Sequence Alignment Overview



Principal component analysis

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

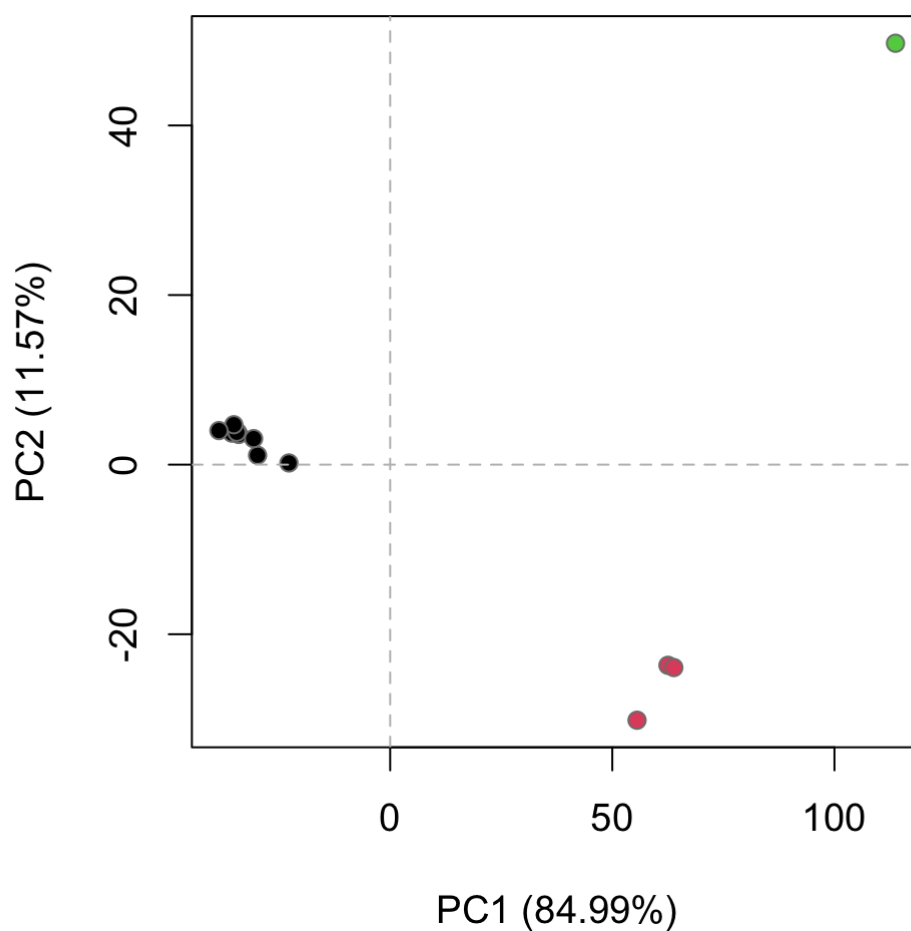


```
# Calculate RMSD
rd <- rmsd(pdbbs)
```

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)
```



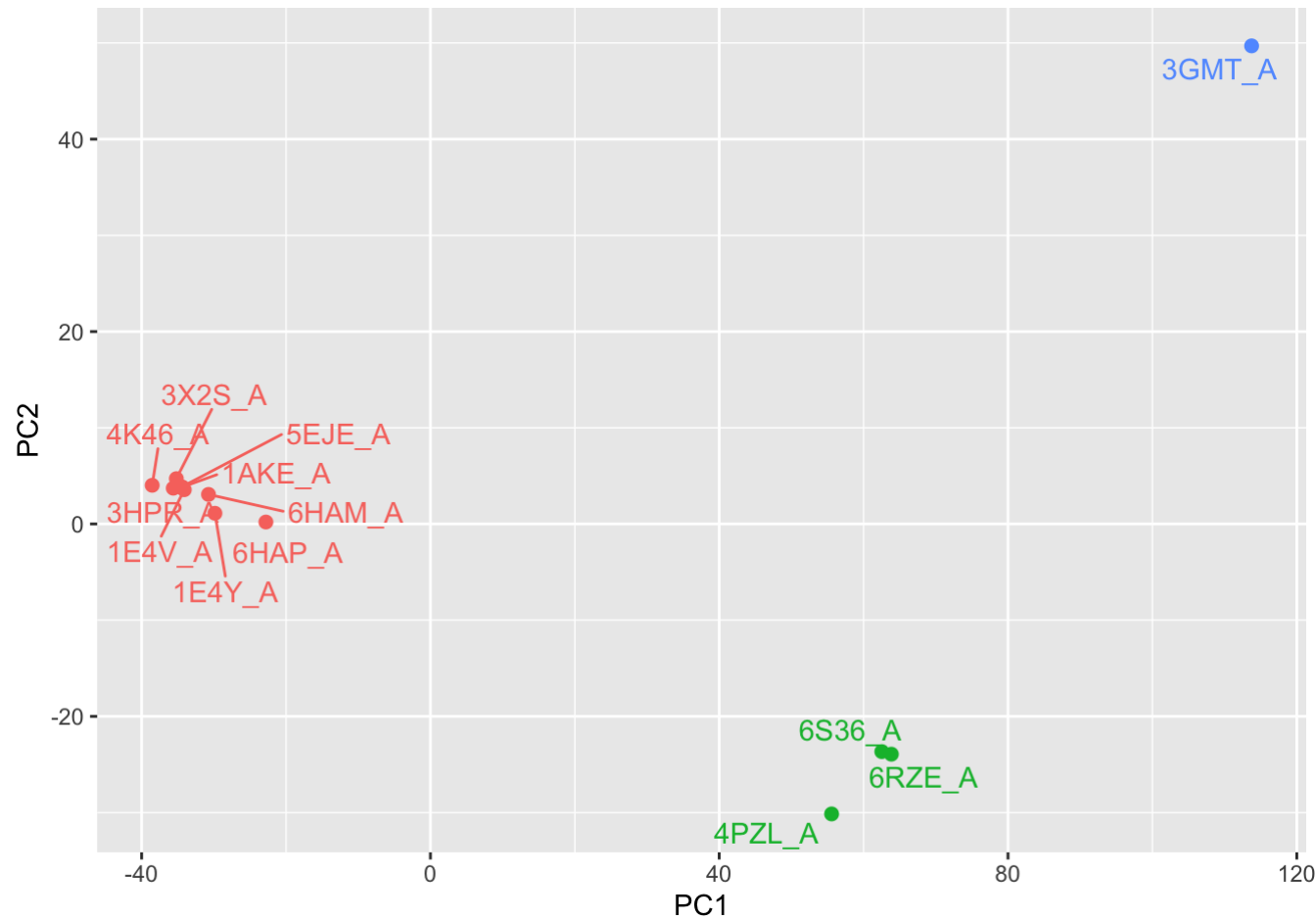
5. Optional further visualization

```
pc1 <- mktrj(pc.xray, pc=1, file="pc_1.pdb")
```

```
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 [optional]

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

Details of Scheduled Calculation:

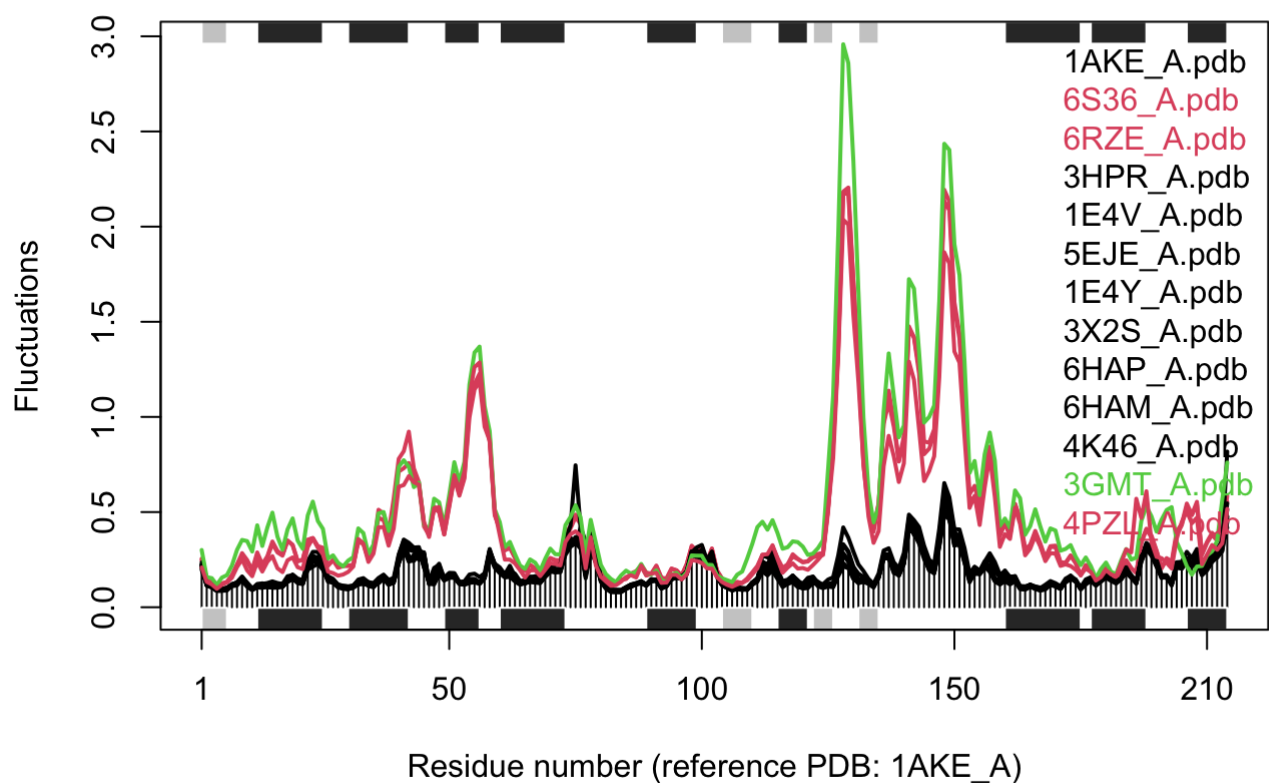
- ... 13 input structures
- ... storing 606 eigenvectors for each structure
- ... dimension of x\$U.subspace: (612x606x13)
- ... coordinate superposition prior to NM calculation
- ... aligned eigenvectors (gap containing positions removed)
- ... estimated memory usage of final 'eNMA' object: 36.9 Mb

| | | |
|-------|--|-----|
| | | 0% |
| | | |
| | | |
| ===== | | 8% |
| | | |
| ===== | | 15% |
| | | |
| ===== | | 23% |



```
plot(modes, pdb, 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?

black and colored look similar in shape but the colored lines have bigger magnitude.