

# Class.10

Samuel Fisher (A18131929)

## Structural Bioinformatics (Part 1)

```
#PDB Statistics
```

```
library(readr)
library(dplyr)
```

```
Attaching package: 'dplyr'
```

```
The following objects are masked from 'package:stats':
```

```
filter, lag
```

```
The following objects are masked from 'package:base':
```

```
intersect, setdiff, setequal, union
```

```
pdb <- read_csv("pdb_stats.csv")
```

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

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

```
Delimiter: ","
```

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

```
dbl (8): X-ray, EM, NMR, Integrative, Multiple methods, Neutron, Other, 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.
```

```
pdb
```

```
# 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>
```

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

```
total_structures <- sum(pdb$Total)
xray_pct <- 100 * sum(pdb$`X-ray`) / total_structures
em_pct <- 100 * sum(pdb$EM) / total_structures
xray_pct
```

```
[1] 80.95077
```

```
em_pct
```

```
[1] 12.83843
```

80.95% of structures in the PDB are solved by X-ray, 12.84% are solved by Electron Microscopy.

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

```
protein_total <- pdb |>
  filter(grepl("Protein", `Molecular Type`)) |>
  summarise(sum_total = sum(Total)) |>
  pull(sum_total)
protein_prop <- protein_total / total_structures
protein_prop
```

```
[1] 0.979118
```

97.91% of structures in the PDB are protein.

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?

There are 1,108 HIV-1 protease structures currently in the PDB.

### **Visualizing the HIV-1 protease structure**

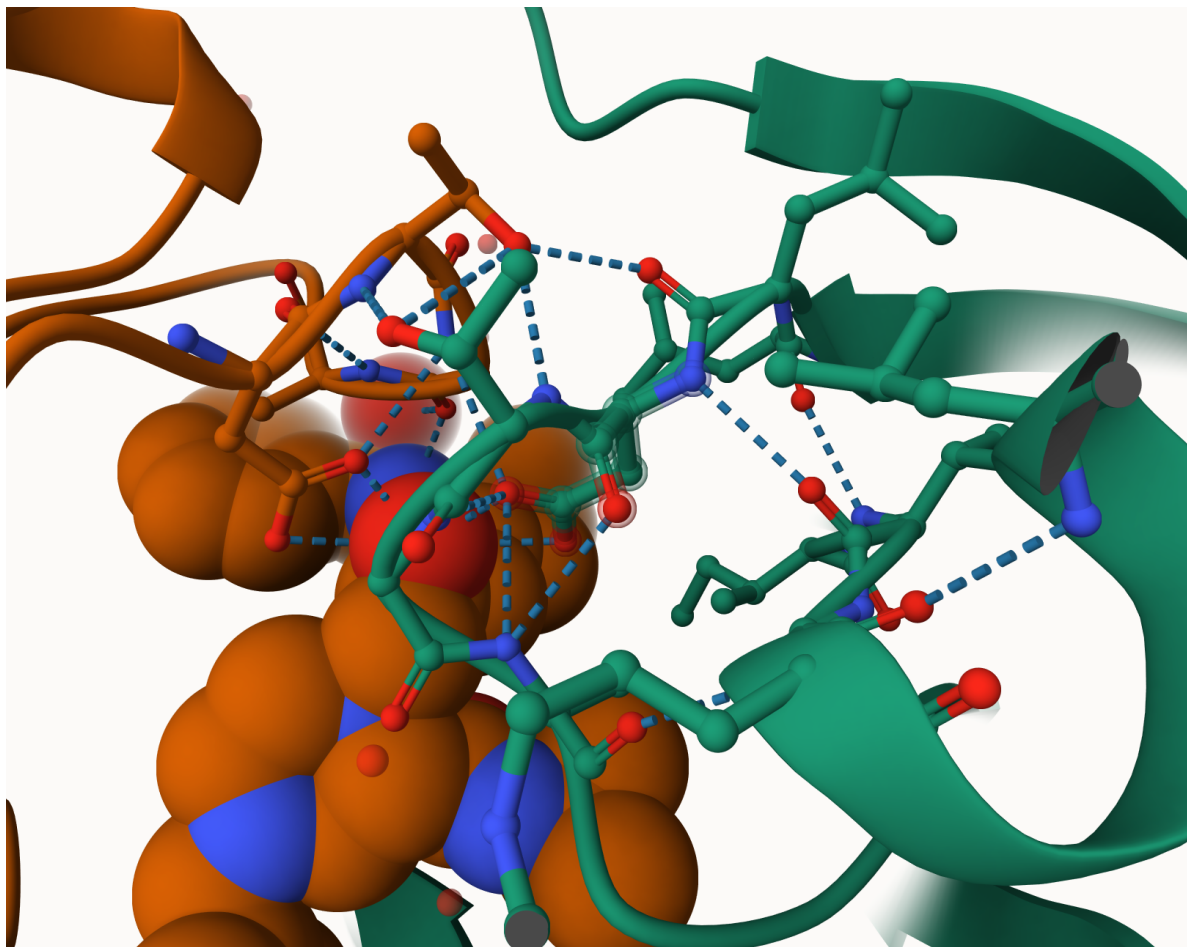
Q4. Water molecules normally have 3 atoms. Why do we see just one atom per water molecule in this structure?

We only see one atom per water molecule because in X-ray crystal structures like 1HSG, hydrogen atoms are usually not resolved. Only the oxygen atom of each water molecule is modeled and shown.

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

HOH 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 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?





The image shows the two HIV-1 protease chains, the bound ligand in the active site, the ASP 25 catalytic residues from each chain, and the conserved water molecule.

#### Section 4: Introduction to Bio3D in R

```
library(bio3d)
```

```
library(bio3dview)
```

```
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:
```

```
PQITLWQRPLVTIKIGGQLKEALLDTGADDTVLEEMSLPGRWKPKMIGGIGGFIVRQYD  
QILIEICGHKAIGTVLVGPTPVNIIGRLLTQIGCTLNFPQITLWQRPLVTIKIGGQLKE  
ALLDTGADDTVLEEMSLPGRWKPKMIGGIGGFIVRQYDQILIEICGHKAIGTVLVGPTP  
VNIIGRLLTQIGCTLNF
```

```
+ 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

## Quick PDB visualization

Unable to generate 3D Protein Models

```
#library(bio3dview)
#library(NGLViewer)

#view.pdb(pdb) |>
  #setSpin()
```

```
#sele <- atom.select(pdb, resno=25)

# and highlight them in spacefill representation
#view.pdb(pdb, cols=c("navy","teal"),
  #highlight = sele,
  #highlight.style = "spacefill") |>
# setRock()
```

## 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:

```
MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMMLRAAVKSGSELGKQAKDIMDAGKLV
DELVIALVKERIAQEDCRNGFLDGFPRTPQADAMKEAGINVDYVLEFDVPELIVDKI
VGRRVHAPSGRVYHVKFNPPKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQMTAPLIG
```

YYSKEAEAGNTKYAKVDGTPVAEVRADLEKILG

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

```
library(bio3d)  
adk <- read.pdb("6s36")
```

Note: Accessing on-line PDB file

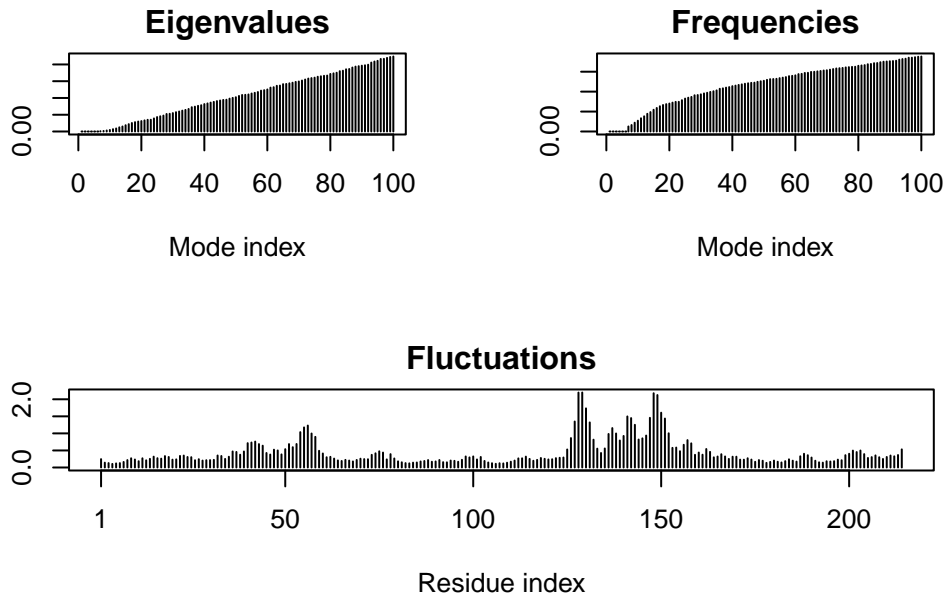
```
Warning in get.pdb(file, path = tempdir(), verbose = FALSE):  
/var/folders/pd/_h44c3lx0kj3kcgk7hp6bx4m0000gn/T//Rtmpx01u48/6s36.pdb exists.  
Skipping download
```

PDB has ALT records, taking A only, rm.alt=TRUE

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

```
Building Hessian...      Done in 0.013 seconds.  
Diagonalizing Hessian... Done in 0.055 seconds.
```

```
plot(m)
```



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

```
Building Hessian...      Done in 0.013 seconds.  
Diagonalizing Hessian... Done in 0.049 seconds.
```

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

```
#view.nma(m, pdb=adk)
```

## Section 5: 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?:

Bio3dview

Q12. True or False? Functions from the pak 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  MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMLRAAVKSGSELGKQAKDIMDAGKLV  
      1      .      .      .      .      .      .      60  
  
      61      .      .      .      .      .      .      120  
pdb|1AKE|A  DELVIALVKERIAQEDCRNGFLLDGFPRPTIPQADAMKEAGINVDYVLEFDVPDELIVDR  
      61      .      .      .      .      .      .      120
```

```

      121      .      .      .      .      .      .      180
pdb|1AKE|A  VGRRVHAPSGRVYHVKFNPPKVEGKDDVTGEELTRKDDQEETVRKRLVEYHQMTAPLIG
      121      .      .      .      .      .      .      180

      181      .      .      .      214
pdb|1AKE|A  YYSKEAEAGNTKYAKVDGTKPVAEVRADLEKILG
      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

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

Download related 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

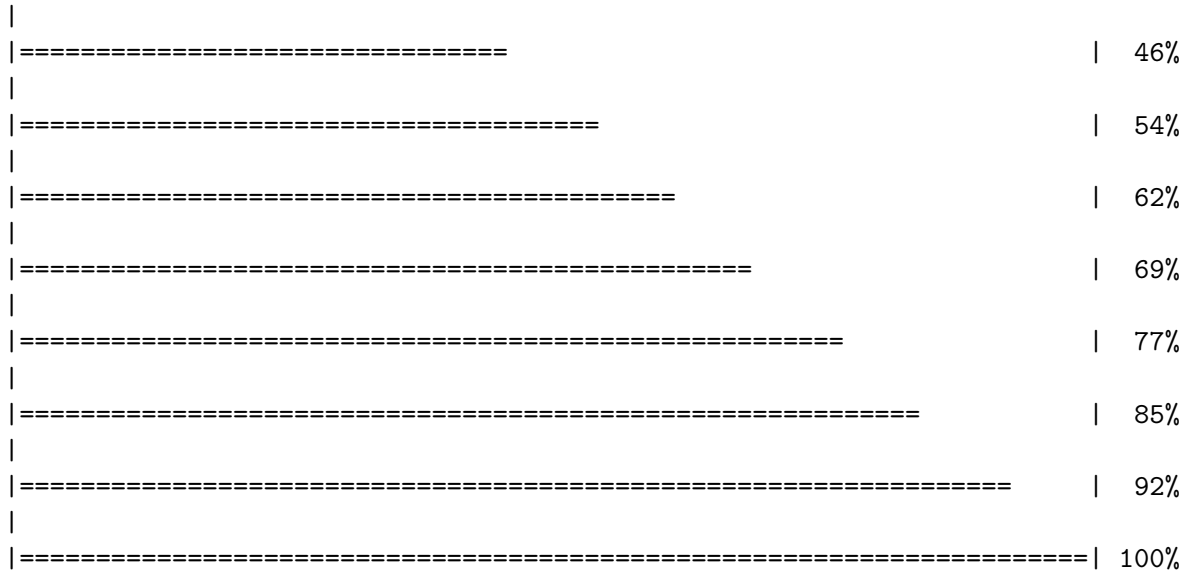
Warning in get.pdb(hits\$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):  
pdbs/6HAM.pdb.gz exists. Skipping download

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%
=====		23%
=====		31%
=====		38%



Align related PDBs

```
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: 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
```

```
#library(bio3dview)
#view.pdbc(pdbc)
```

## Annotate

Vector containing PDB database codes

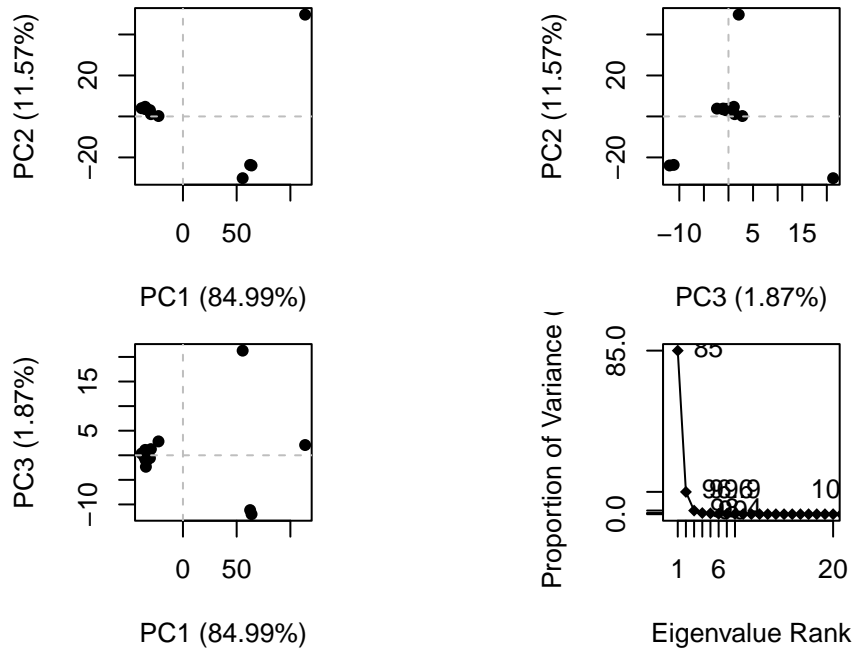
```
ids <- basename.pdbc(pdbc$id)
anno <- pdbc.annotate(ids)
unique(anno$source)
```

```
[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"
```

Perform PCA

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



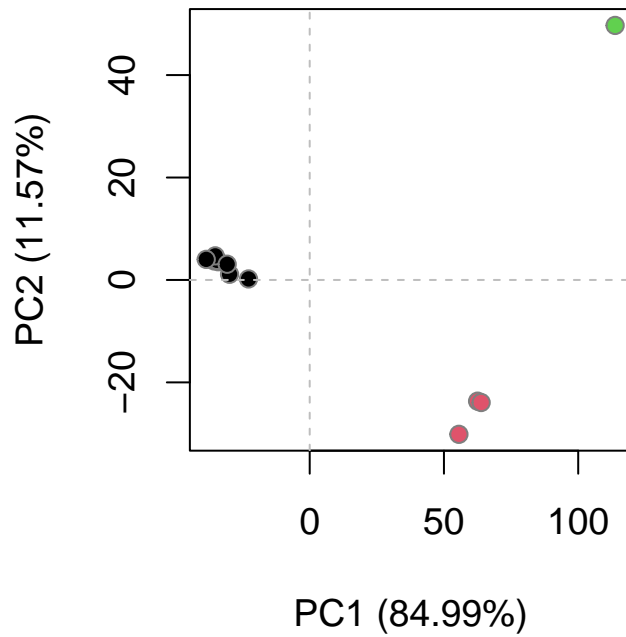
Calculate RMSD

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

Warning in rmsd(pdfs): 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 <- mktrj(pc.xray, pc=1, file="pc_1.pdb")
```

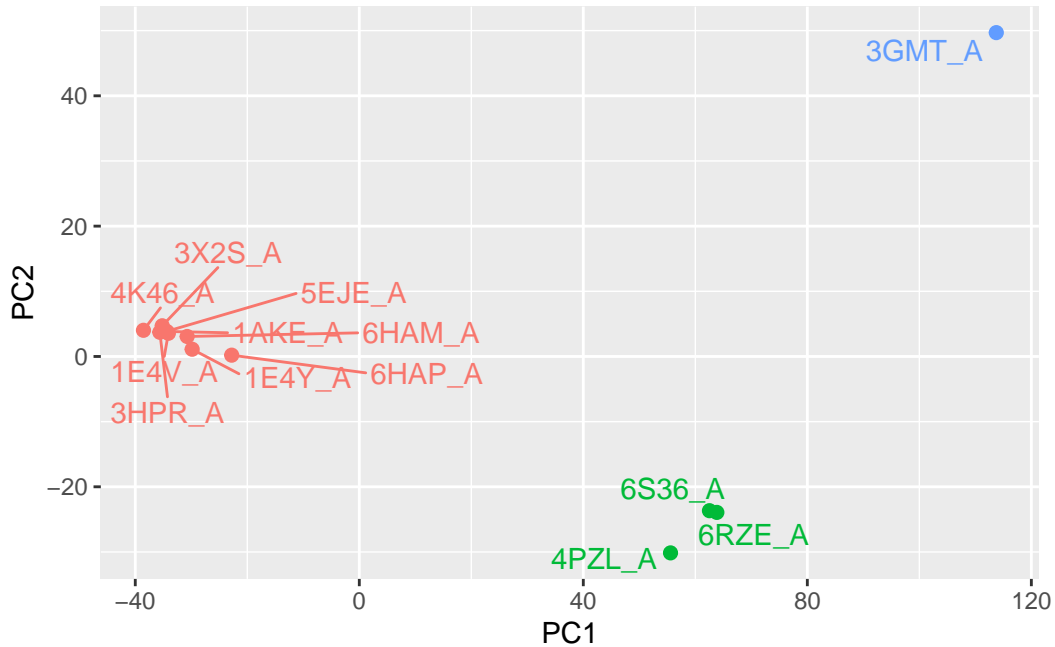
lotting results with ggplot2

```
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
```



## Comparative protein structure analysis with PCA

We start with a database id "lake\_A"

```
library(bio3d)

id <- "lake_A"
aa <- get.seq(id)
```

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

Fetching... Please wait. Done.

```
#blast <- blast.pdb(aa)
```

Have a peak:

```
#head(blast$hit.tbl)
```

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

Peak at our “top hits”

```
head(hits$pdb.id)
```

```
[1] "1AKE_A" "6S36_A" "6RZE_A" "3HPR_A" "1E4V_A" "5EJE_A"
```

Now we can download these “top hits” these will all be ADK structures in the PDB database.

```
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
```

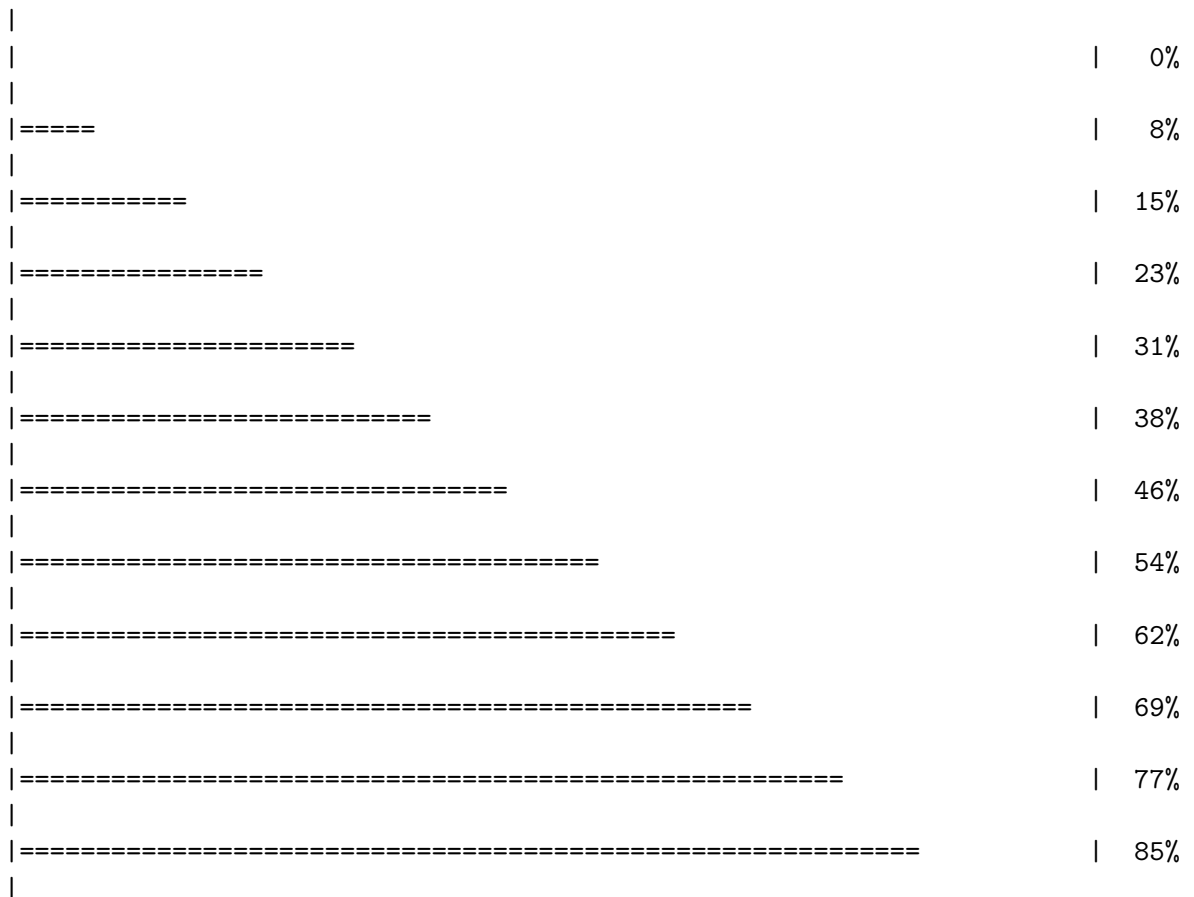
Warning in get.pdb(hits\$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):  
pdbs/6HAP.pdb.gz exists. Skipping download

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



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

We need one package from BioConductor. To set this up we need to first install a package called **BiocManager** from CRAN Now we can use the `install()` function from this package like this: `BiocManager::install("msa")`

```
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

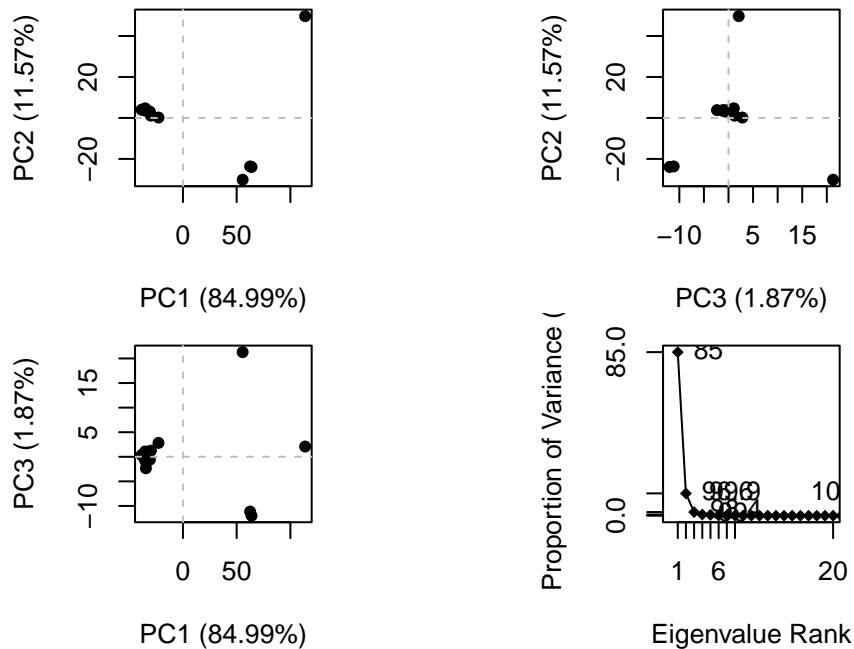
```

Let's have a little peak at our structures after "fitting" or superposing:

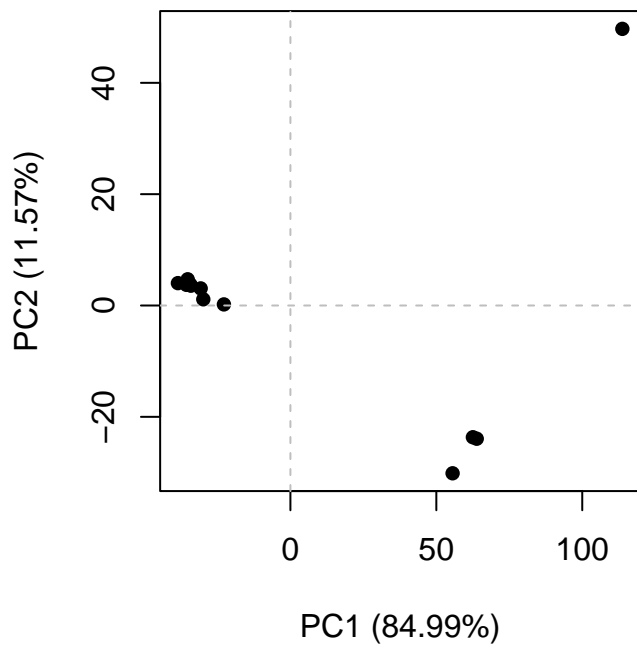
```
#view.pdbs(pdbs,colorScheme="residue")
```

We can run functions like `rmsd()`, `rmsd()` and the best `pca()`

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



```
plot(pc.xray, 1:2)
```



Finally, let's make a movie of the major "motion" or structural difference in the dataset - we call this a "trajectory".

```
mktrj(pc.xray, file = "results.pdb")
```