

Lab 09

AUTHOR

Nichelle Camden

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

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

169,916 out of 196,979 total, about 86%.

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

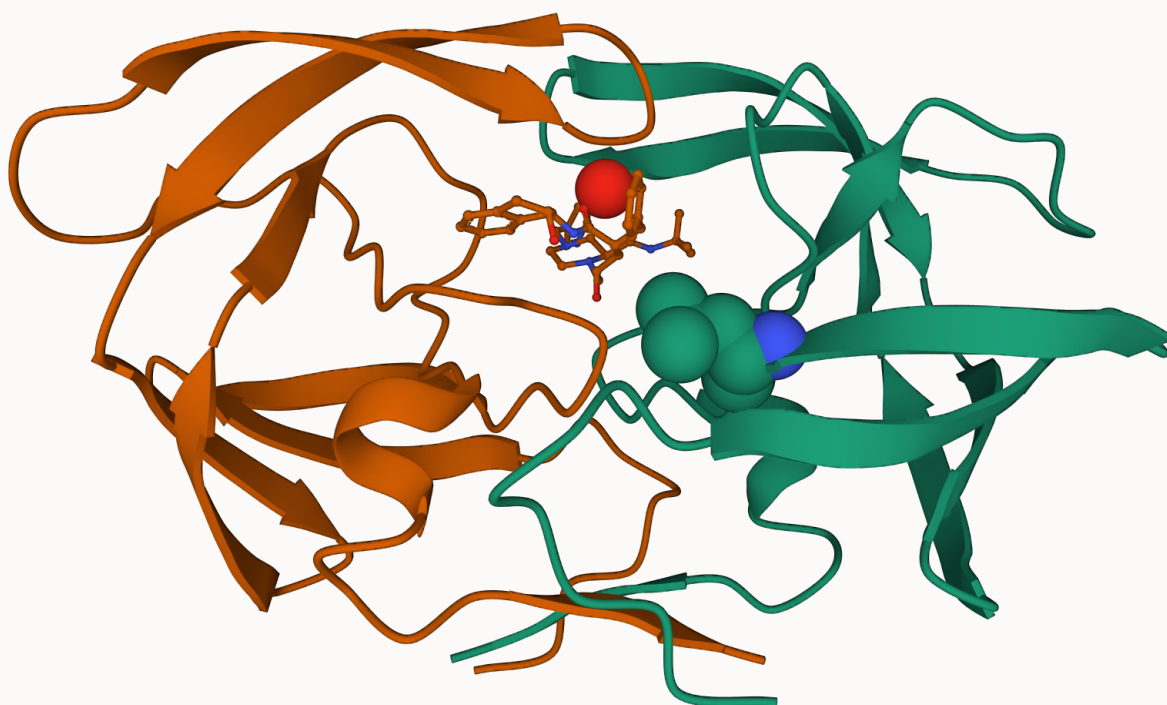
most of them! 171,351 out of 196,979 total, about 87%

#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?

22,183

2. Visualizing the HIV-1 protease structure

Viewing PDB structures with Mol*



Reading and working with structures in R

The `bio3d` package for structural bioinformatics has a lot of features for reading and working with biomolecular sequences and structures.

```
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:

DOTTI WARDI VTTKTGGDI KEAI IDTGADDTVI FEESI DGRWKPKMTGGTGGETKVRQVD

```

FQITLWQRFLVTRIGGQERALEDTGADDTVEELPSLGRWKFNFIGGIGGIKVRQTD
QILIEICGHKAIGTVLVGPTPVNIIGRNLLTQIGCTLNFPQITLWQRPLVTIKIGGQLKE
ALLDTGADDTVLEEMSLPGRWKPKMIGGIGGFIKVRQYDQILIEICGHKAIGTVLVGPTP
VNIIGRNLLTQIGCTLNF

```

```

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

```

```
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
TDELVIALVKERIAQEDCRNGFLLDGFPRITPQADAMKEAGINVDYVLEFDVPDELIVDKI
VGRRVHAPSGRVYHVKFNPVKVEGKDDVTGEELTRKDDQEETVRKRLVEYHQMTAPLIG
YYSKEAEAGNTKYAKVDGTPVAEVRADLEKILG

```

```

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

```

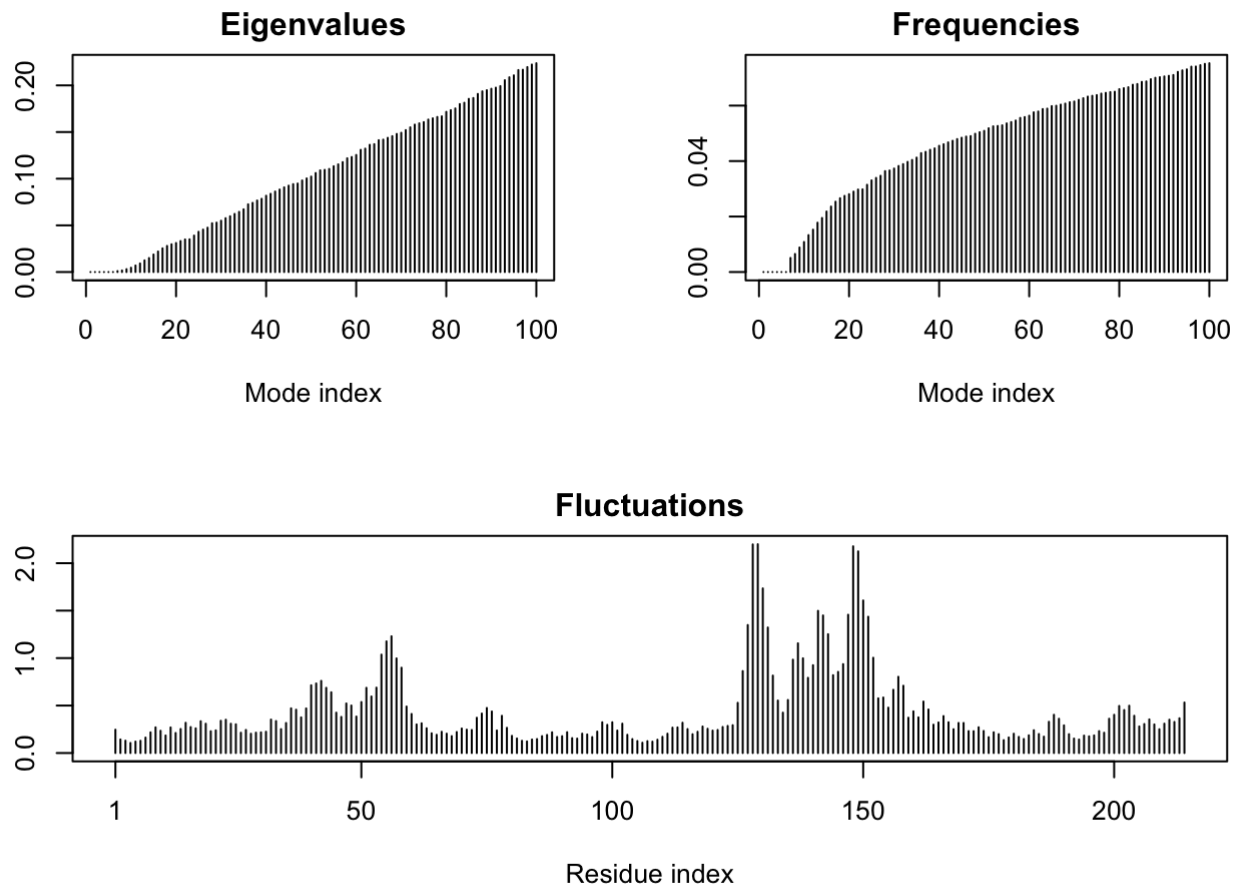
Normal mode analysis (NMA) is a bioinformatics method for predicting functional motions. It will show us the parts of the protein that are "flexible" (i.e. the most dynamic).

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

```
Building Hessian... Done in 0.084 seconds.
```

```
Diagonalizing Hessian... Done in 0.634 seconds.
```

```
plot(m)
```



Make a "movie" of this thing moving.

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

Comparative analysis of all ADK structures

First, we get the sequence of ADK and use this to search the PDB database.

```
aa <- get.seq("1ake_a")
```

Warning in get.seq("1ake_a"): Removing existing file: seqs.fasta

Fetching... Please wait. Done.

```
aa
```

```

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

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

```

```

pdb|1AKE|A  VGGRRVHAPSGRVYHVKNPPKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQMTAPLIG
121          .          .          .          .          .          120
121          .          .          .          .          .          180
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
```

Run a blast search of the PDB database using this sequence

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

Searching ... please wait (updates every 5 seconds) RID = P02CFPWJ016

.....

Reporting 98 hits

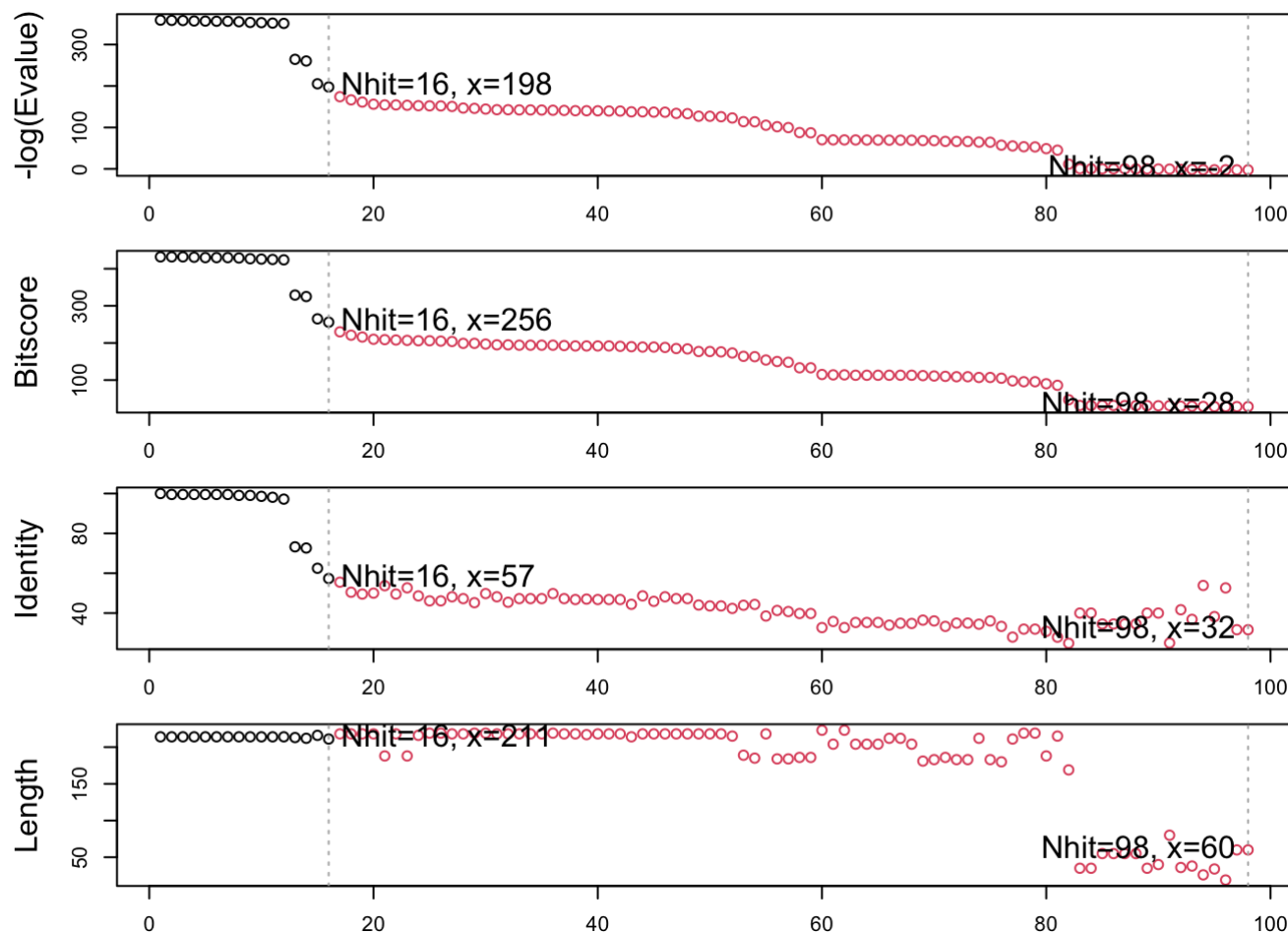
```
#blast
```

examine the results

```
hits <- plot(blast)
```

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

```
* Chosen cutoff value of: 197
  Yielding Nhits: 16
```



download all the ADK structures in the PDB database

```
hits$pdb.id
```

```
[1] "1AKE_A" "4X8M_A" "6S36_A" "6RZE_A" "4X8H_A" "3HPR_A" "1E4V_A" "5EJE_A"
[9] "1E4Y_A" "3X2S_A" "6HAP_A" "6HAM_A" "4K46_A" "4NP6_A" "3GMT_A" "4PZL_A"
```

```
#pdb.annotate(hits$pdb.id)
```

```
# Download related PDB files
```

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

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

```
Warning in get.pdb(hits$pdb.id, path = "pds", split = TRUE, gzip = TRUE): pds/
4X8M.pdb.gz exists. Skipping download
```

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

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

Warning in get.pdb(hits\$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE): pdbs/
4X8H.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/
4NP6.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%
====	6%
=====	12%
=====	19%
=====	25%
=====	31%
=====	38%



Viewing all these structures looks like a mess. We need to try something else...

We will align and superimpose these structures.

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

Reading PDB files:

```
pdbbs/split_chain/1AKE_A.pdb
pdbbs/split_chain/4X8M_A.pdb
pdbbs/split_chain/6S36_A.pdb
pdbbs/split_chain/6RZE_A.pdb
pdbbs/split_chain/4X8H_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/4NP6_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

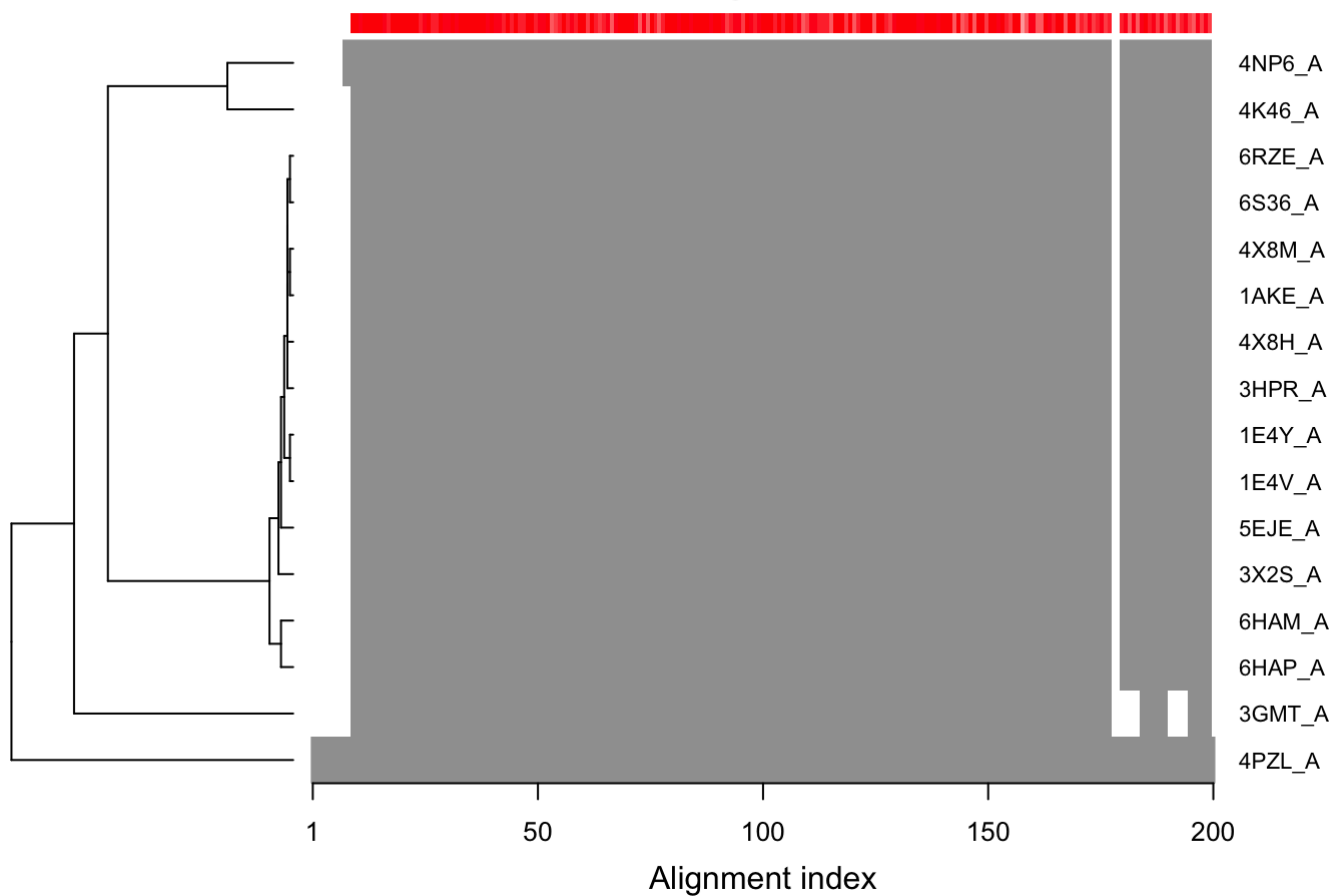
```
pdb/seq: 1 name: pdbs/split_chain/1AKE_A.pdb
PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 2 name: pdbs/split_chain/4X8M_A.pdb
pdb/seq: 3 name: pdbs/split_chain/6S36_A.pdb
PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 4 name: pdbs/split_chain/6RZE_A.pdb
PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 5 name: pdbs/split_chain/4X8H_A.pdb
pdb/seq: 6 name: pdbs/split_chain/3HPR_A.pdb
PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 7 name: pdbs/split_chain/1E4V_A.pdb
pdb/seq: 8 name: pdbs/split_chain/5EJE_A.pdb
PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 9 name: pdbs/split_chain/1E4Y_A.pdb
pdb/seq: 10 name: pdbs/split_chain/3X2S_A.pdb
pdb/seq: 11 name: pdbs/split_chain/6HAP_A.pdb
pdb/seq: 12 name: pdbs/split_chain/6HAM_A.pdb
PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 13 name: pdbs/split_chain/4K46_A.pdb
PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 14 name: pdbs/split_chain/4NP6_A.pdb
pdb/seq: 15 name: pdbs/split_chain/3GMT_A.pdb
pdb/seq: 16 name: pdbs/split_chain/4PZL_A.pdb
```

```
#pdbs
```

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

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

Sequence Alignment Overview

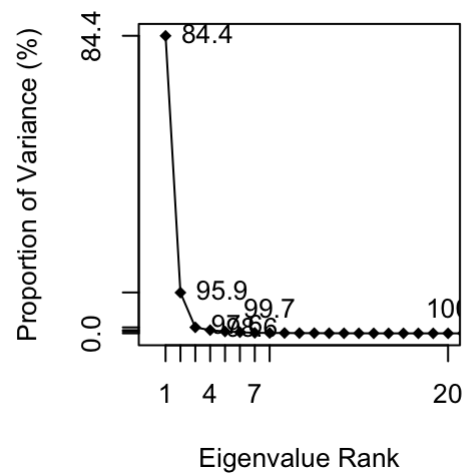
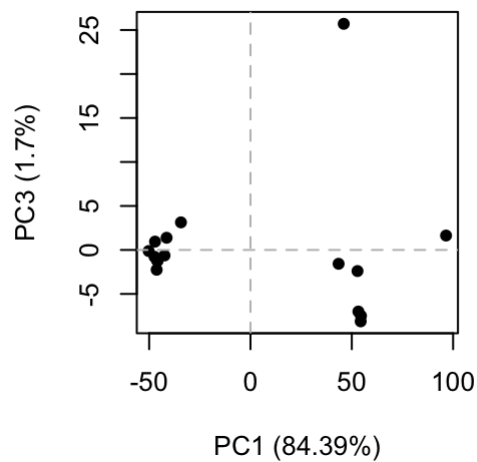
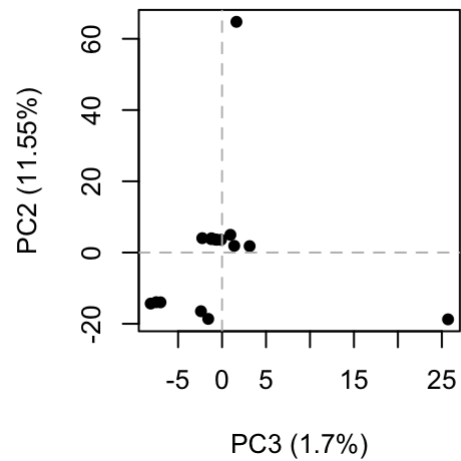
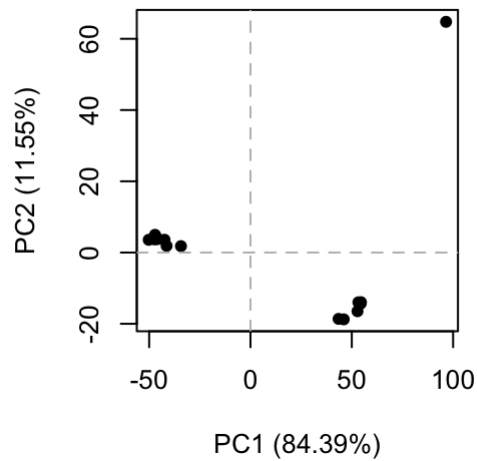


PCA to the rescue...

```
pc.xray <- pca(pdb)
```

and plot my results

```
plot(pc.xray)
```



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

Warning in rmsd(pdb): 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)
```

