# Class_9_structural_bioinformatics

Max Strul

10-26-2022

## Table of contents

## Getting started viewing pdb files

You can use multiple viewing tools such as Chimerax or the free online website molstart `https://molstar.org/viewer/`

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

```
total_xray <- 150417/171351
total_xray
```

```
[1] 0.8778297
```

```
total_em <- 8586/171351
total_em
```

```
[1] 0.05010767
```

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

```
171351/(171351+10459+10919+4037+191+22)
```

```
[1] 0.8698948
```

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?

4707

# bio3d and its function in loading PDB files

To load and interact with pdb files into R we need a package `bio3d`

```
library(bio3d)
```

```
example_pdb <-read.pdb("1hsg")
```

```
 Note: Accessing on-line PDB file
```

```
example_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
```

```
  head(example_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>
```

Looking at other proteins: Adenylate Kinase

```
  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:
   MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMLRAAVKSGSELGKQAKDIMDAGKLVT
   DELVIALVKERIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFDVPDELIVDKI
   VGRRVHAPSGRVYHVKFNPPKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQMTAPLIG
   YYSKEAEAGNTKYAKVDGTKPVAEVRADLEKILG

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

We can use important information here such as the "(residues/Calpha atoms#: 214)" which tells us there are 214 amino acids.

We can look at all other attributes.

```
attributes(adk)
```

```
$names
[1] "atom"   "xyz"    "seqres" "helix"  "sheet"  "calpha" "remark" "call"

$class
[1] "pdb" "sse"
```

We can analyze the functional motions of structure

## Predicting functional motions of a single structure

The first thing we will do is plot the normal mode analysis. This is a structural bioinformatics method to determine and predict the protein flexibility and potential functional motions.
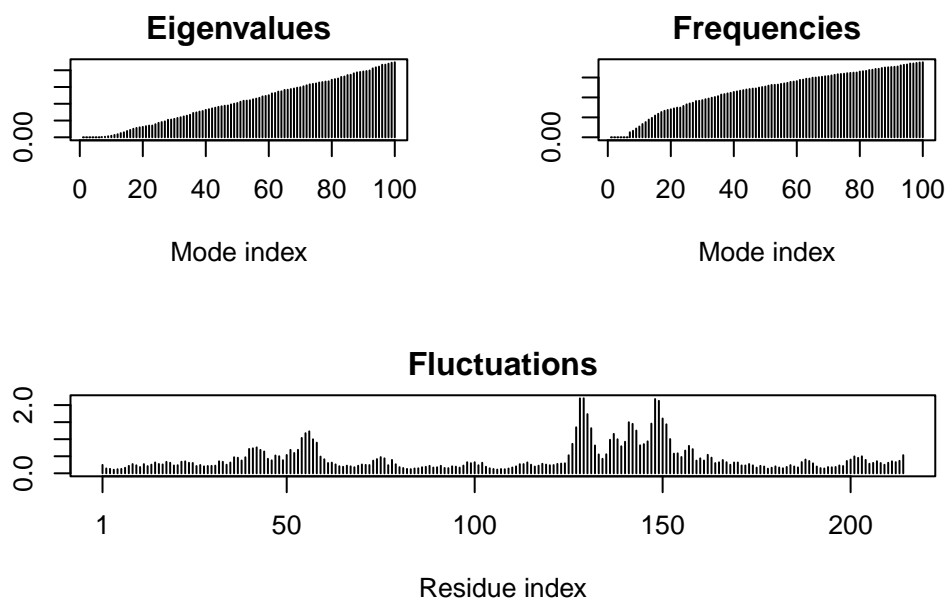
```
# Perform flexiblity prediction
m <- nma(adk)
```

```
Building Hessian...        Done in 0.083 seconds.
Diagonalizing Hessian...   Done in 0.257 seconds.
```

```r
plot(m)
```

**Eigenvalues**

Mode index

**Frequencies**

Mode index

**Fluctuations**

Residue index

Now we want to use this to make a movie and plot the trajectory of certain residues

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

This function produces a file that has the multiple states and can be opened in molstar to view a movie of the enzyme different conformation changes as calculated by the normal mode analysis (which is performed via a force field calculation modeling atomic bonds as springs)

Now what we want to do is obtain all homolog of this protein and see what their conformation changes are like, and which ones are active and which ones are inactive, etc.

## Using R to create a repository of homologs in order to do a comparative analysis across structures

We first need to load in all the predicted similar protein structures:

First we want to get our initial protein of interest, in this case it is ADK.

5

We will then `get.seq()` on this protein, then use this sequence to search a protein database (PDB database) with the function `blast.pdb()`, then we will get a list of all homologs and get their .pdb files with `get.pdb()`

```
#aminoacidseq <- get.seq("1ake_A")
#blast_results <- blast.pdb(aminoacidseq) #commented out but usually you would run this
#we can plot a summary plot
#hits <- plot(blast_results)
#this graph has the Y axis in units of -log(Evalue)

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

We can see that there are some Nhits = # where there are some amount that are good enough based on their E value

```
hits$pdb.id
```

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

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

Now we want to obtain our zipped/folder package of our 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%
  |
```

```
  |==============================                               |  46%
  |
  |=====================================                        |  54%
  |
  |============================================                 |  62%
  |
  |==================================================           |  69%
  |
  |=========================================================    |  77%
  |
  |============================================================ |  85%
  |
  |===============================================================|  92%
  |
  |================================================================| 100%
```

Now we want to overlap all of them:

```
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: 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
```

```r
#pdbs
```
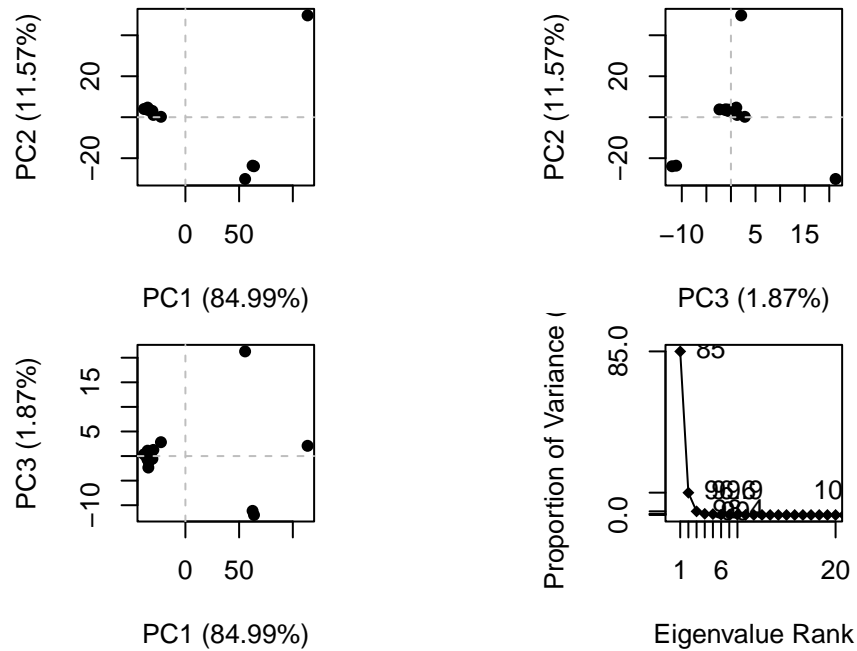
```r
ids <- basename.pdb(pdbs$id)
```

```r
#plot(pdbs, labels=ids)
```

(sequencealignment.png)

Then we want to start our comparative analysis with a principal comparative analysis

## Principal component analysis
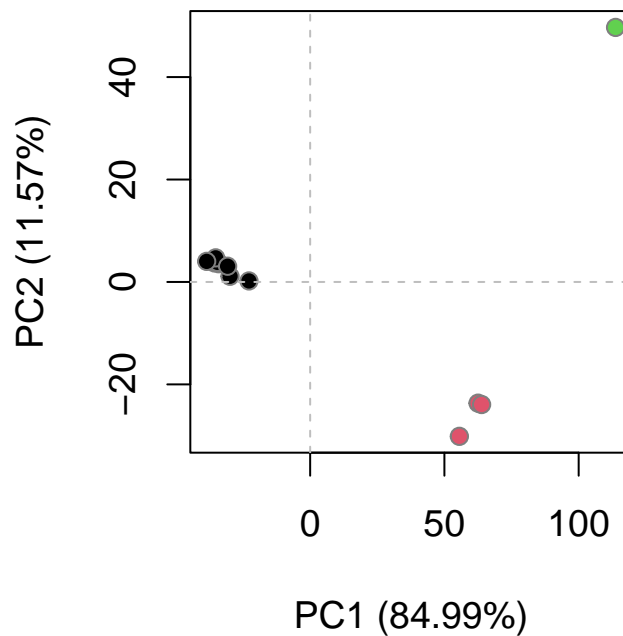
```r
pc.xray <- pca(pdbs)
plot(pc.xray)
```

You can do calculate RMSD to do structure-based clustering

```
rd <- rmsd(pdbs)
```

Warning in rmsd(pdbs): 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)
```

Here we see there are clearly two different groups

When we plot this and when we perform mktrj we are essentially creating dots between the two different structural groups (marked here in red and black)

When we compare inhibitor some are more inhibited than others, is there any way that this plot determines this? Short answer; no, not in this initial plotting at least. It is possible that if you were to characterize the inhibition of certain structures and make that another variable mask you could theoretically perform a PCA to determine a separating function that figures out how you could develop and make a trajectory between low inhibition and more inhibition
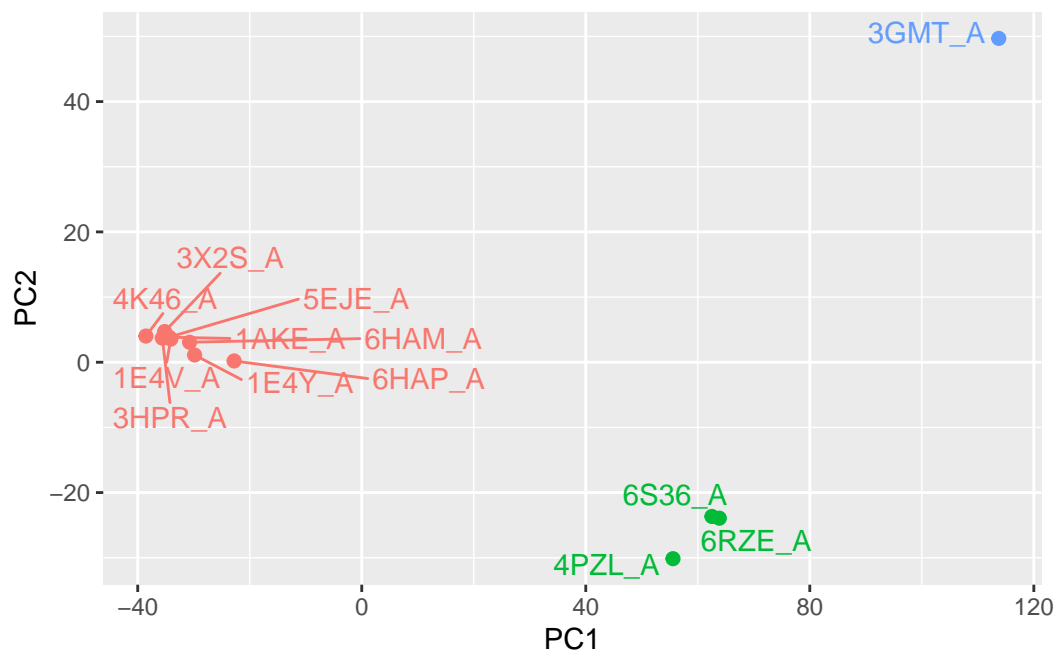
```
pc1 <- mktrj(pc.xray, pc=1, file="pc_1.pdb")


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
```



```
modes <- nma(pdbs)
```

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

```r
plot(modes, pdbs, col=grps.rd)
```

Extracting SSE from pdbs$sse attribute