

# Class 9: Structural Bioinformatics 1

Kira

## What is the PDB anyway?

The main database of biomolecular structures is called the PDB and is available at [www.rcsb.org](http://www.rcsb.org).

Let's begin by seeing what is in this database:

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

	X.ray	EM	NMR	Multiple.methods	Neutron	Other
Protein (only)	152,809	9,421	12,117	191	72	32
Protein/Oligosaccharide	9,008	1,654	32	7	1	0
Protein/NA	8,061	2,944	281	6	0	0
Nucleic acid (only)	2,602	77	1,433	12	2	1
Other	163	9	31	0	0	0
Oligosaccharide (only)	11	0	6	1	0	4
Total						
Protein (only)	174,642					
Protein/Oligosaccharide	10,702					
Protein/NA	11,292					
Nucleic acid (only)	4,127					
Other	203					
Oligosaccharide (only)	22					

## How can we change the values in the table from strings to numbers?

```
as.numeric(gsub(",", "", pdbstats$X.ray))
```

```
[1] 152809    9008    8061    2602    163     11
```

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

```
# Getting rid of the commas so we can do math
n.xray <- sum(as.numeric(gsub(",", "", pdbstats$X.ray)))
n.em <- sum(as.numeric(gsub(",", "", pdbstats$EM)))
n.total <- sum(as.numeric(gsub(",", "", pdbstats$Total)))
```

```
p.xray <- (n.xray)/(n.total) * 100
p.em <- (n.em)/(n.total) * 100
```

```
# and to 2 s.f
round(p.xray, 2)
```

```
[1] 85.9
```

```
round(p.em, 2)
```

```
[1] 7.02
```

The percentage of X-ray structures is 85.9% and the percentage of EM structures is 7.02% (in the current PDB database).

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

```
# Looking at totals of all rows, the first output will correspond to protein only
as.numeric(gsub(",", "", pdbstats$Total))/n.total * 100
```

```
[1] 86.89175473  5.32469600  5.61824587  2.05335642  0.10100105  0.01094593
```

The proportion of structures in the PDB that are only protein is 86.89%.

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?

It is not straightforward to find all HIV-1 protease structures using plain text searching on the database.

## 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?

Hydrogen atoms are so small that they often cannot be seen in the X-ray crystallography structure of the protein (i.e the hydrogen atom is smaller than the resolution of the structure).

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 is a critically conserved water molecule that is important to the binding site.

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.

## Introduction to Bio3D in R

We will use the `bio3d` package for this section:

```
library(bio3d)
```

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

Note: Accessing on-line PDB file

```
pdb
```

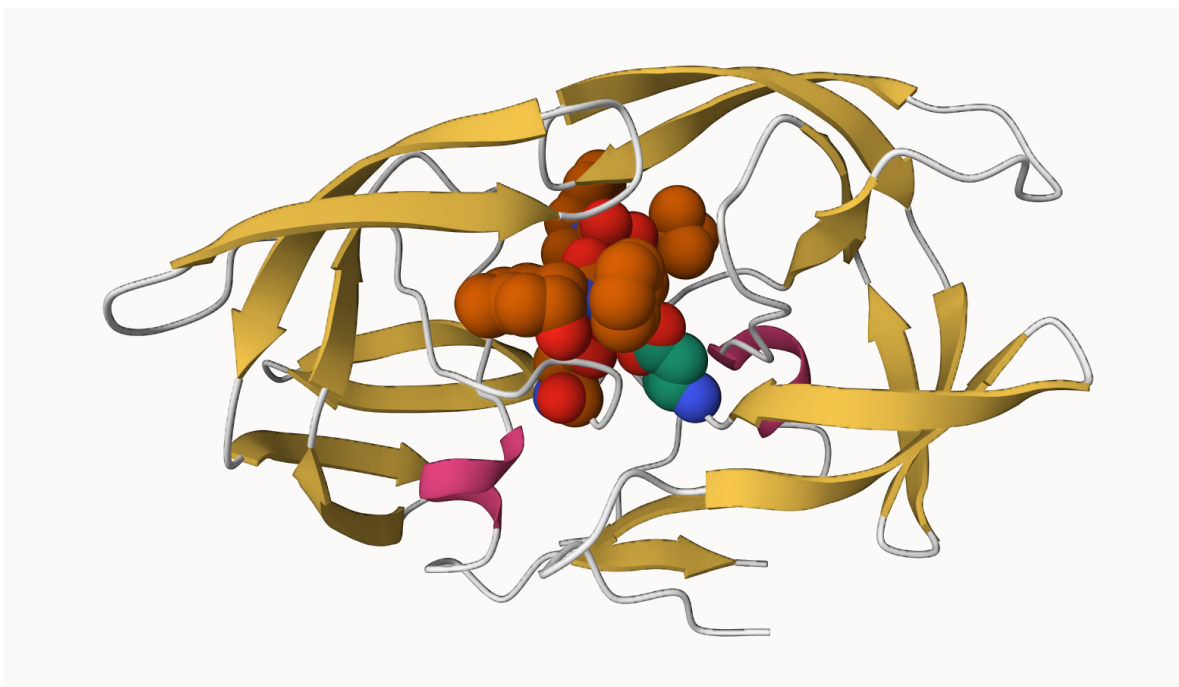


Figure 1: Asp25 residues and HOH 308 are spacefilled in addition to ligand.

```
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?

There are 198 amino acid residues.

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

HOH.

Q9: How many protein chains are in this structure?

There are 2 protein chains in this structure.

```
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>

What is the first residue 3 letter code and 1 letter code?:

```
pdb$atom$resid[1]
```

```
[1] "PRO"
```

```
aa321(pdb$atom$resid[1])
```

```
[1] "P"
```

## Predicting functional motions of a single structure

Let's read a new PDB structure of Adenlyate Kinase (PDB code: 6s36) and perform normal mode analysis.

```
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  
VGRRVHAPSGRVYHVKFNPVKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQMTAPLIG  
YYSKEAEAGNTKYAKVDGTPVAEVRADLEKILG
```

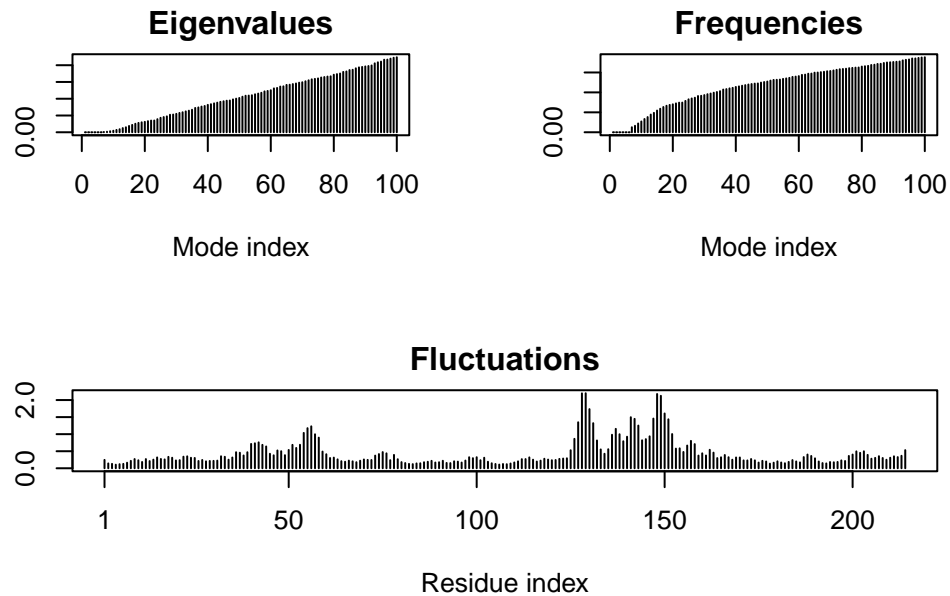
```
+ attr: atom, xyz, seqres, helix, sheet,  
      calpha, remark, call
```

Normal mode analysis (NMA) is a structural bioinformatics method to predict protein flexibility and potential functional motions (a.k.a. conformational changes).

```
# Now we can perform normal mode analysis  
m <- nma(adk)
```

```
Building Hessian...      Done in 0.075 seconds.  
Diagonalizing Hessian... Done in 0.656 seconds.
```

```
plot(m)
```



Peaks in the fluctuations plot show areas where the protein is most flexible.

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

## Comparative structure analysis of Adenylate Kinase

Today we are continuing where we left off on Tuesday, building towards completing the loop from biomolecular structural data to our new analysis methods like PCA and clustering.

We will begin by getting a single protein sequence for a protein family of interest.

```
# Install packages in the R console NOT your Rmd/Quarto file  
  
#install.packages("bio3d")  
#install.packages("devtools")
```

```
#install.packages("BiocManager")

#BiocManager::install("msa")
#devtools::install_bitbucket("Grantlab/bio3d-view")
```

Q10. Which of the packages above is found only on BioConductor and not CRAN?

`msa()` is found only on BioConductor and not CRAN.

Q11. Which of the above packages is not found on BioConductor or CRAN?:

`biod3d.view()` is not found on BioConductor or CRAN.

Q12. True or False? Functions from the `devtools` package can be used to install packages from GitHub and BitBucket?

This is true, we can use `devtools()` to install packages from both GitHub and BitBucket.

```
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  DELVIALVKERIAQEDCRNGFLLDGFPRPTIPQADAMKEAGINVDYVLEFDVPDELIVDRI
      61      .      .      .      .      .      .      120

      121      .      .      .      .      .      .      180
pdb|1AKE|A  VGRRVHAPSGRVYHVKNPPKVEGKDDVTGEELTRKDDQEETVRKRLVEYHQMTAPLIG
      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?

There are 214 amino acids in this sequence.

Now we can use this sequence as a query to BLAST search the PDB to find similar sequences and structures.

```
# Blast or hmmer search, will comment out second line to prevent subsequent runs when re-run
# b <- blast.pdb(aa)
```

I could save and load my blast results next time so I don't need to run the search every time.

```
# saveRDS(b,file="blast_results.RDS")
```

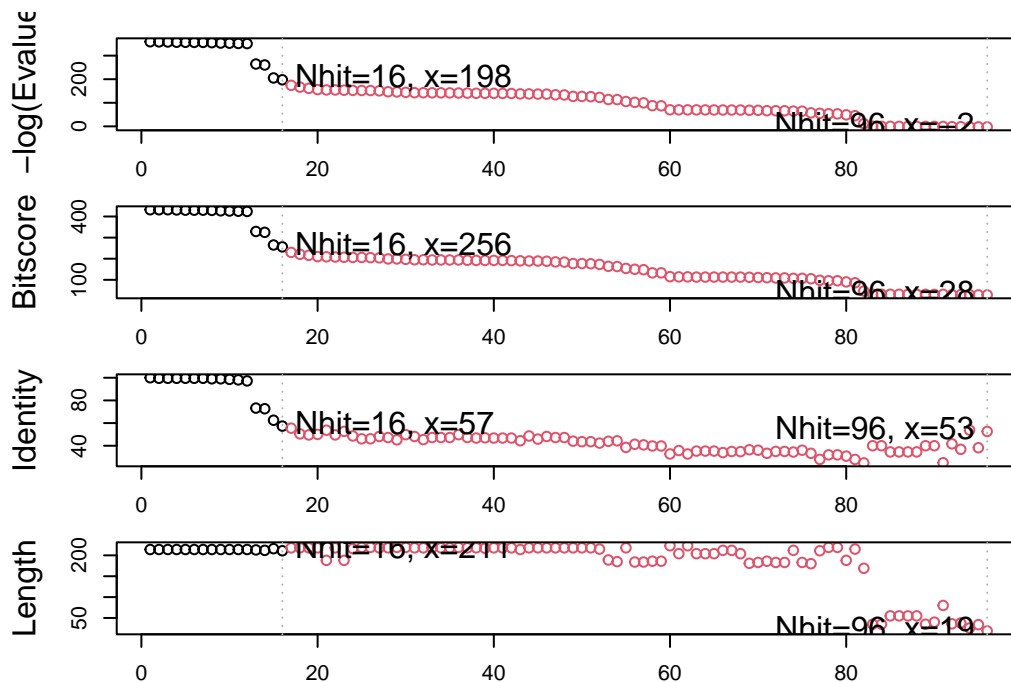
```
# If I wanted to see the results again
b <- readRDS("blast_results.RDS")
```

We now want to generate a summary of our BLAST results using a plot.

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

```
* Possible cutoff values: 197 -3
    Yielding Nhits:      16 96
```

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



```
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"
```

```
# Download related PDB files, additional arguments will make our data a little tidier
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/4X8M.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/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

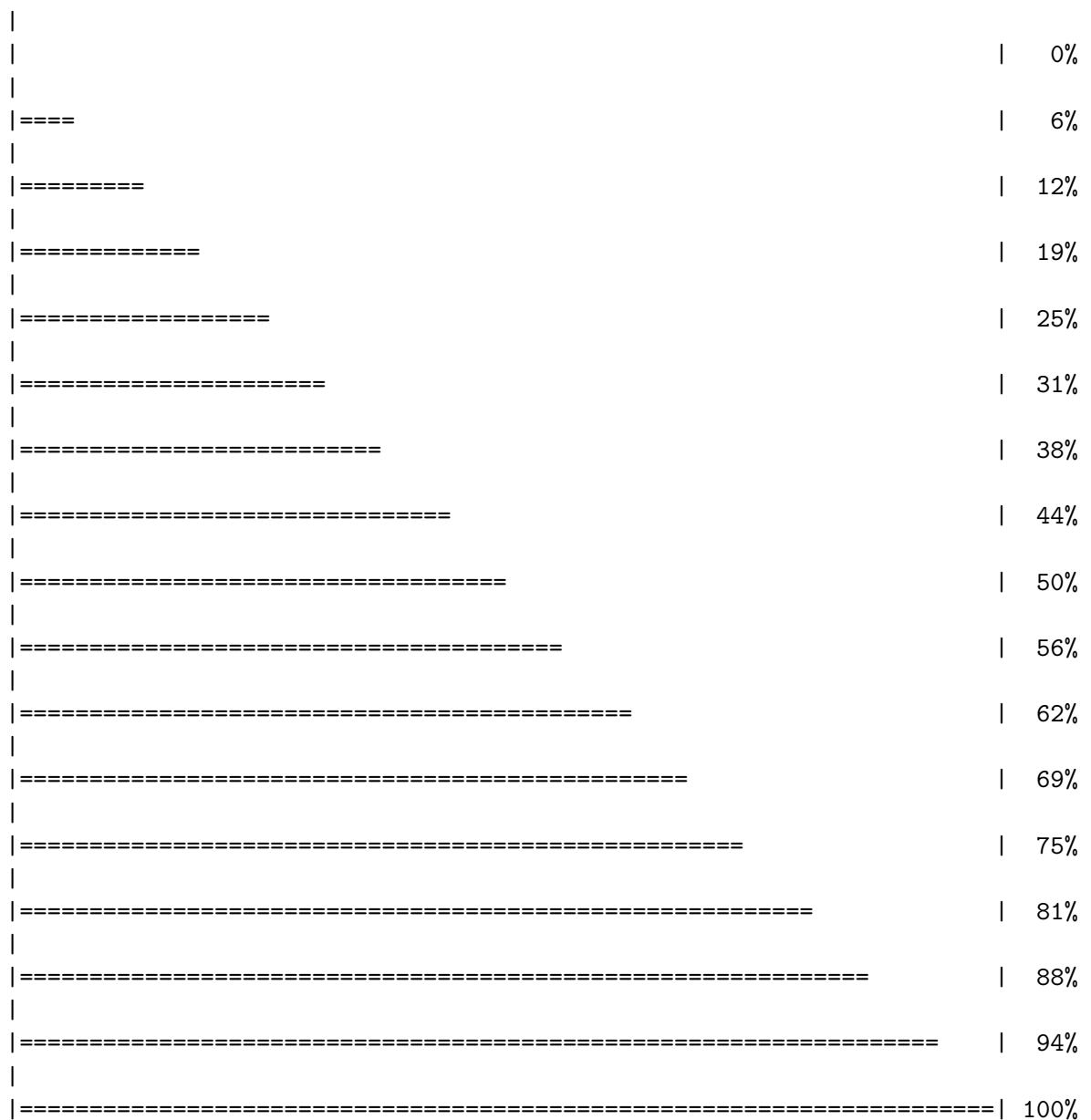
Warning in get.pdb(hits\$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):  
pdbs/6HAM.pdb.gz exists. Skipping download

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



Next, we are going to align and superimpose all these structures.

## Align and superpose structures

```
# Align related PDBs. fit = TRUE will superimpose them and "msa" will allow us to perform
pdbs <- pdbaln(files, fit = TRUE, exefile="msa")
```

Reading PDB files:

```
pdbs/split_chain/1AKE_A.pdb
pdbs/split_chain/4X8M_A.pdb
pdbs/split_chain/6S36_A.pdb
pdbs/split_chain/6RZE_A.pdb
pdbs/split_chain/4X8H_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/4NP6_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/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: pdbc/split_chain/3HPR_A.pdb
            PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 7   name: pdbc/split_chain/1E4V_A.pdb
pdb/seq: 8   name: pdbc/split_chain/5EJE_A.pdb
            PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 9   name: pdbc/split_chain/1E4Y_A.pdb
pdb/seq: 10  name: pdbc/split_chain/3X2S_A.pdb
pdb/seq: 11  name: pdbc/split_chain/6HAP_A.pdb
pdb/seq: 12  name: pdbc/split_chain/6HAM_A.pdb
            PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 13  name: pdbc/split_chain/4K46_A.pdb
            PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 14  name: pdbc/split_chain/4NP6_A.pdb
pdb/seq: 15  name: pdbc/split_chain/3GMT_A.pdb
pdb/seq: 16  name: pdbc/split_chain/4PZL_A.pdb

```

## pdbc

```

[Truncated_Name:1] 1AKE_A.pdb      1          .          .          .          40
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:2] 4X8M_A.pdb      1          .          .          .          40
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:3] 6S36_A.pdb      1          .          .          .          40
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:4] 6RZE_A.pdb      1          .          .          .          40
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:5] 4X8H_A.pdb      1          .          .          .          40
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:6] 3HPR_A.pdb      1          .          .          .          40
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:7] 1E4V_A.pdb      1          .          .          .          40
-----MRIILLGAPVAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:8] 5EJE_A.pdb      1          .          .          .          40
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:9] 1E4Y_A.pdb      1          .          .          .          40
-----MRIILLGALVAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:10] 3X2S_A.pdb      1          .          .          .          40
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:11] 6HAP_A.pdb      1          .          .          .          40
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:12] 6HAM_A.pdb      1          .          .          .          40
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:13] 4K46_A.pdb      1          .          .          .          40
-----MRIILLGAPGAGKGTQAQFIMAKFGIPQIS
[Truncated_Name:14] 4NP6_A.pdb      1          .          .          .          40
-----NAMRIILLGAPGAGKGTQAQFIMEKFGIPQIS
[Truncated_Name:15] 3GMT_A.pdb      1          .          .          .          40
-----MRLILLGAPGAGKGTQANFIKEKFGIPQIS
[Truncated_Name:16] 4PZL_A.pdb      1          .          .          .          40
-----TENLYFQSNAMRIILLGAPGAGKGTQAKIIEQYNIAHIS
                **~*****  *****  *  *~ *  **
1          .          .          .          40

41          .          .          .          80
TGDMLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVKE
[Truncated_Name:1] 1AKE_A.pdb      41          .          .          .          80
TGDMLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVKE
[Truncated_Name:2] 4X8M_A.pdb      41          .          .          .          80
TGDMLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVKE

```

[Truncated_Name:3] 6S36_A.pdb	TGDMRLRAAVKSGSELGKQAKDIMDAGKLVTDDELVIALVKE
[Truncated_Name:4] 6RZE_A.pdb	TGDMRLRAAVKSGSELGKQAKDIMDAGKLVTDDELVIALVKE
[Truncated_Name:5] 4X8H_A.pdb	TGDMRLRAAVKSGSELGKQAKDIMDAGKLVTDDELVIALVKE
[Truncated_Name:6] 3HPR_A.pdb	TGDMRLRAAVKSGSELGKQAKDIMDAGKLVTDDELVIALVKE
[Truncated_Name:7] 1E4V_A.pdb	TGDMRLRAAVKSGSELGKQAKDIMDAGKLVTDDELVIALVKE
[Truncated_Name:8] 5EJE_A.pdb	TGDMRLRAAVKSGSELGKQAKDIMDACKLVTDDELVIALVKE
[Truncated_Name:9] 1E4Y_A.pdb	TGDMRLRAAVKSGSELGKQAKDIMDAGKLVTDDELVIALVKE
[Truncated_Name:10] 3X2S_A.pdb	TGDMRLRAAVKSGSELGKQAKDIMDCGKLVTDDELVIALVKE
[Truncated_Name:11] 6HAP_A.pdb	TGDMRLRAAVKSGSELGKQAKDIMDAGKLVTDDELVIALVRE
[Truncated_Name:12] 6HAM_A.pdb	TGDMRLRAAIKSGSELGKQAKDIMDAGKLVTDDEIIIALVKE
[Truncated_Name:13] 4K46_A.pdb	TGDMRLRAAIKAGTELGKQAKSVIDAGQLVSDDIILGLVKE
[Truncated_Name:14] 4NP6_A.pdb	TGDMRLRAAIKAGTELGKQAKAVIDAGQLVSDDIILGLIKE
[Truncated_Name:15] 3GMT_A.pdb	TGDMRLRAAVKAGTPLGVEAKTYMDEGKLVPSLIIGLVKE
[Truncated_Name:16] 4PZL_A.pdb	TGDMIRETIKSGSALGQELKKVLDAGELVSDEFIIVKIVKD
	****~* ~* *~ ** * ~* ** * ~ ~~~~~
	41 . . . 80
	81 . . . 120
[Truncated_Name:1] 1AKE_A.pdb	RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:2] 4X8M_A.pdb	RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:3] 6S36_A.pdb	RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:4] 6RZE_A.pdb	RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:5] 4X8H_A.pdb	RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:6] 3HPR_A.pdb	RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:7] 1E4V_A.pdb	RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:8] 5EJE_A.pdb	RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:9] 1E4Y_A.pdb	RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:10] 3X2S_A.pdb	RIAQEDSRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:11] 6HAP_A.pdb	RICQEDSRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:12] 6HAM_A.pdb	RICQEDSRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:13] 4K46_A.pdb	RIAQDDCAKGFLDGFPR TIPQADGLKEVG VVDYVIEFD
[Truncated_Name:14] 4NP6_A.pdb	RIAQADCEKGFLLDGFPR TIPQADGLKEMGINVDYVIEFD
[Truncated_Name:15] 3GMT_A.pdb	RLKEADCANGYLF DGFPR TIPQADAMKEAGVAIDYVLEID
[Truncated_Name:16] 4PZL_A.pdb	RISKNDCNNGFLLDGVPR TIPQAQELDKLG VNIIDYIVEVD
	*~ * *~* ** ***** ** ^ *~ ~***~* *
	81 . . . 120
	121 . . . 160
[Truncated_Name:1] 1AKE_A.pdb	VPDELIVDRIVGRRVHAPSGRVYHV KFNPPKVEGKDDVTG
[Truncated_Name:2] 4X8M_A.pdb	VPDELIVDRIVGRRVHAPSGRVYHV KFNPPKVEGKDDVTG
[Truncated_Name:3] 6S36_A.pdb	VPDELIVDKIVGRRVHAPSGRVYHV KFNPPKVEGKDDVTG
[Truncated_Name:4] 6RZE_A.pdb	VPDELIVDAIVGRRVHAPSGRVYHV KFNPPKVEGKDDVTG
[Truncated_Name:5] 4X8H_A.pdb	VPDELIVDRIVGRRVHAPSGRVYHV KFNPPKVEGKDDVTG





```

[Truncated_Name:9]1E4Y_A.pdb      T--KYAKVDGTPVAEVRADLEKILG-
[Truncated_Name:10]3X2S_A.pdb     T--KYAKVDGTPVAEVRADLEKILG-
[Truncated_Name:11]6HAP_A.pdb     T--KYAKVDGTPVCEVRADLEKILG-
[Truncated_Name:12]6HAM_A.pdb     T--KYAKVDGTPVCEVRADLEKILG-
[Truncated_Name:13]4K46_A.pdb     T--QYLKFDGTKAVAEVSAELEKALA-
[Truncated_Name:14]4NP6_A.pdb     T--QYLKFDGTKQVSEVSADIAKALA-
[Truncated_Name:15]3GMT_A.pdb     E-----NGLKAPA-----YRKISG-
[Truncated_Name:16]4PZL_A.pdb     KIPKYIKINGDQAVEKVSQDIFDQLNK
                                   *
                                   .           .           227
201

```

Call:

```
pdbaln(files = files, fit = TRUE, exefile = "msa")
```

Class:

```
pdbs, fasta
```

Alignment dimensions:

```
16 sequence rows; 227 position columns (204 non-gap, 23 gap)
```

```
+ attr: xyz, resno, b, chain, id, ali, resid, sse, call
```

```

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

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

```

In this plot, grey regions depict aligned residues, while white depict gap regions. The red bar at the top depicts sequence conservation.

We also want to collect annotation for each entry:

```

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] "Vibrio cholerae 01 biovar El Tor str. N16961"

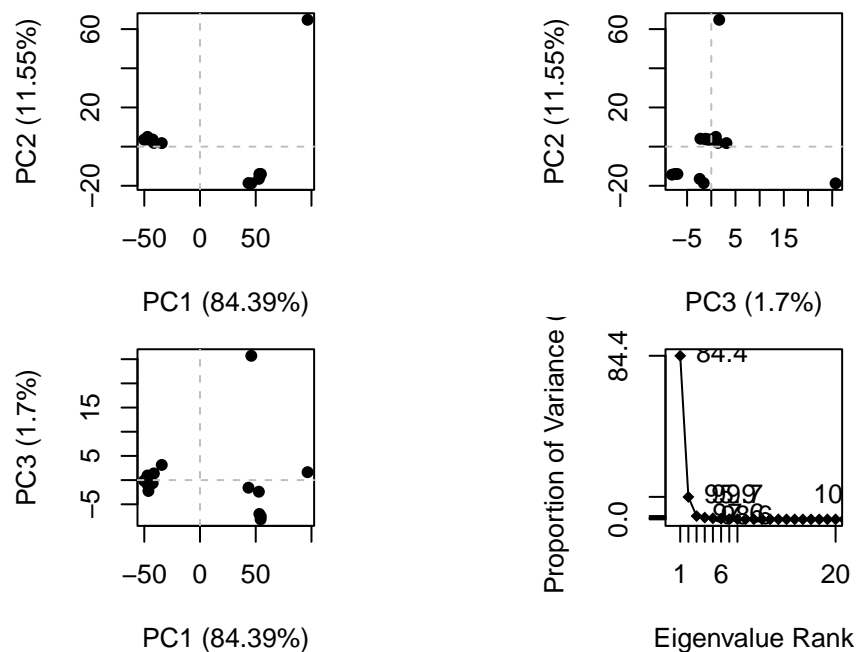
```

```
[7] "Burkholderia pseudomallei 1710b"
[8] "Francisella tularensis subsp. tularensis SCHU S4"
```

## Principal component analysis of Adenylate kinase

Now we can do PCA to try and make sense of this data. We will not use the `prcomp()` function from base R, but the `pca()` function from the `bio3d` package as it is designed to work nicely with biomolecular data.

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



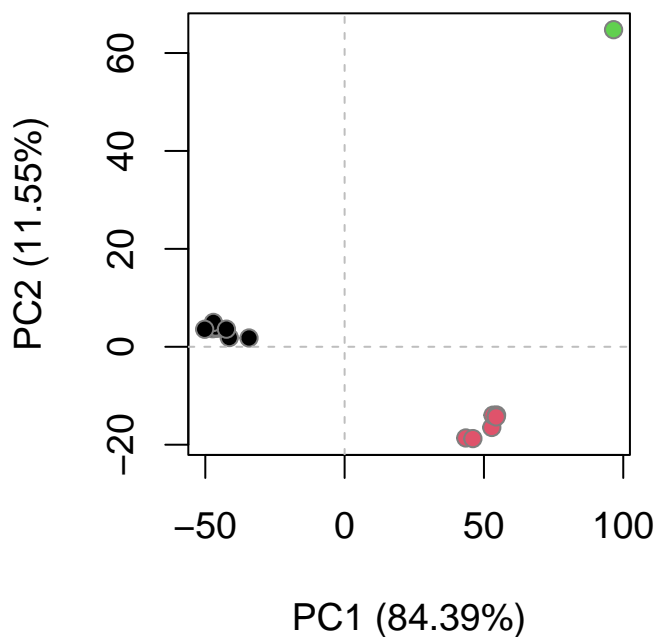
We can now focus in on PC1 vs. PC2 and make a plot that is colored by cluster group.

```
# Calculate RMSD (distance between all elements)
rd <- rmsd(pdbbs)
```

Warning in `rmsd(pdbbs)`: No indices provided, using the 204 non NA positions

```
# Structure-based clustering (making 3 groups)
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)
```



Now, we want to visualize major structural variations:

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

We can now open this trajectory file in Molstar to view a movie of the major differences (i.e. displacements) in the structure set as we move along PC1.

We can also plot this information in `ggplot()` for a better looking, more informational plot.

```
#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

```

