

# Class 9: Structural bioinformatics 1

Derek Chang (PID: 16942232)

The main database for structural data is called PDB (Protein Data Bank). Let's see what it contains:

Data: <https://tinyurl.com/pdbstats24>

```
pdbdb <- read.csv("pdb_stats.csv")
pdbdb
```

	Molecular.Type	X.ray	EM	NMR	Multiple.methods	Neutron	Other
1	Protein (only)	167,192	15,572	12,529	208	77	32
2	Protein/Oligosaccharide	9,639	2,635	34	8	2	0
3	Protein/NA	8,730	4,697	286	7	0	0
4	Nucleic acid (only)	2,869	137	1,507	14	3	1
5	Other	170	10	33	0	0	0
6	Oligosaccharide (only)	11	0	6	1	0	4
	Total						
1		195,610					
2		12,318					
3		13,720					
4		4,531					
5		213					
6		22					

Questions:

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

```
pdbdb$Total
```

```
[1] "195,610" "12,318" "13,720" "4,531" "213" "22"
```

I need to remove the commas and convert to numeric to do math:

```
as.numeric(sub(",", "", pdbdb$Total))
```

```
[1] 195610 12318 13720 4531 213 22
```

I could turn this snippet into a function to fix any future table i read like this:

```
x <- pdbdb$Total  
as.numeric(sub(",", "", x))
```

```
[1] 195610 12318 13720 4531 213 22
```

```
comma2numeric <- function(x) {  
  as.numeric(sub(",", "", x))  
}
```

Test it

```
pdbdbnew<- comma2numeric(pdbdb$X.ray)
```

##Or try a different read/import function:

```
library(readr)  
pdbdb <- read_csv("pdb_stats.csv")
```

Rows: 6 Columns: 8

-- Column specification -----

Delimiter: ","

chr (1): Molecular Type

dbl (3): Multiple methods, Neutron, Other

num (4): X-ray, EM, NMR, 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.

```
sum(pdbdb$Total)
```

```
[1] 226414
```

```
sum(pdbdb$`X-ray`)/sum(pdbdb$Total) * 100
```

```
[1] 83.30359
```

```
sum(pdbdb$EM)/sum(pdbdb$Total) * 100
```

```
[1] 10.18091
```

83.3% of the structures are solved by X-ray and 10.18% of the structures are solved by EM.

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

```
pdbdb
```

```
# A tibble: 6 x 8
```

`Molecular Type` <chr>	`X-ray` <dbl>	EM <dbl>	NMR <dbl>	`Multiple methods` <dbl>	Neutron <dbl>	Other <dbl>	Total <dbl>
1 Protein (only)	167192	15572	12529	208	77	32	195610
2 Protein/Oligosacc~	9639	2635	34	8	2	0	12318
3 Protein/NA	8730	4697	286	7	0	0	13720
4 Nucleic acid (onl~	2869	137	1507	14	3	1	4531
5 Other	170	10	33	0	0	0	213
6 Oligosaccharide (~	11	0	6	1	0	4	22

```
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

```
(sum(pdbdb$Total[1])/sum(pdbdb$Total)) *100
```

```
[1] 86.39483
```

86.39483% of the 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?

In the current PDB there are 226414 HIV-1 Protease Structures.

## 2. Mol\*

Mol\* is a new web-based molecular viewer that we will use to learn the basics of here.

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

We see just one atom per water molecule as this allows for better clarity of the water molecules and allows us to see other things better.

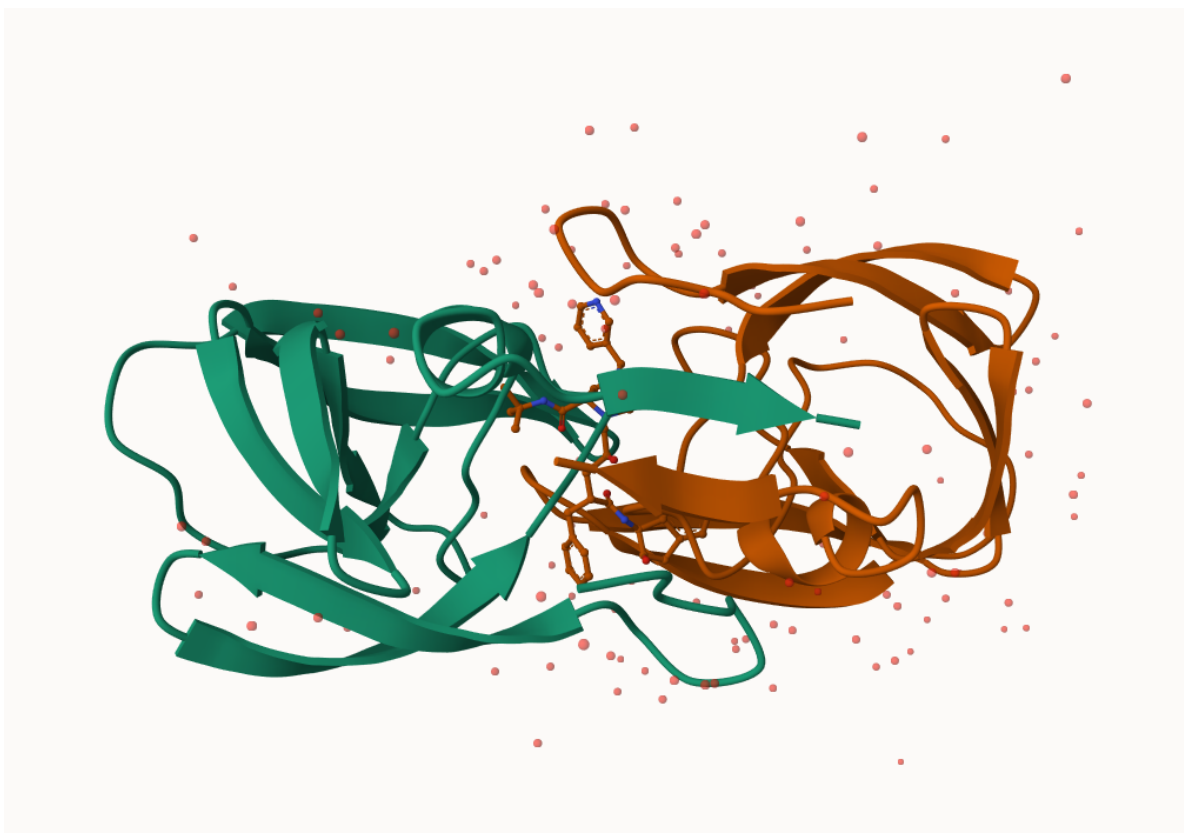


Figure 1: A first image from molstar

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

The conserved water molecule in the binding site seems to be residue number 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.



Figure 2: A second image from molstar showing the important ASP25 amino acid molecules.

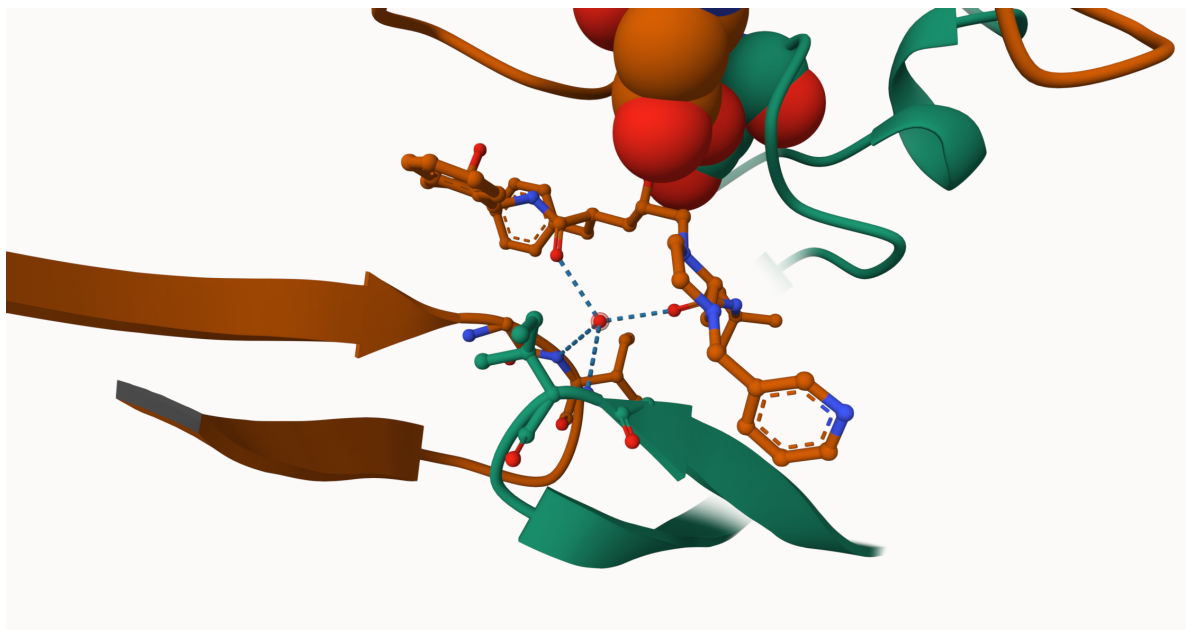


Figure 3: Image showing the two distinct chains of HIV-protease and the critical conserved water residue 308

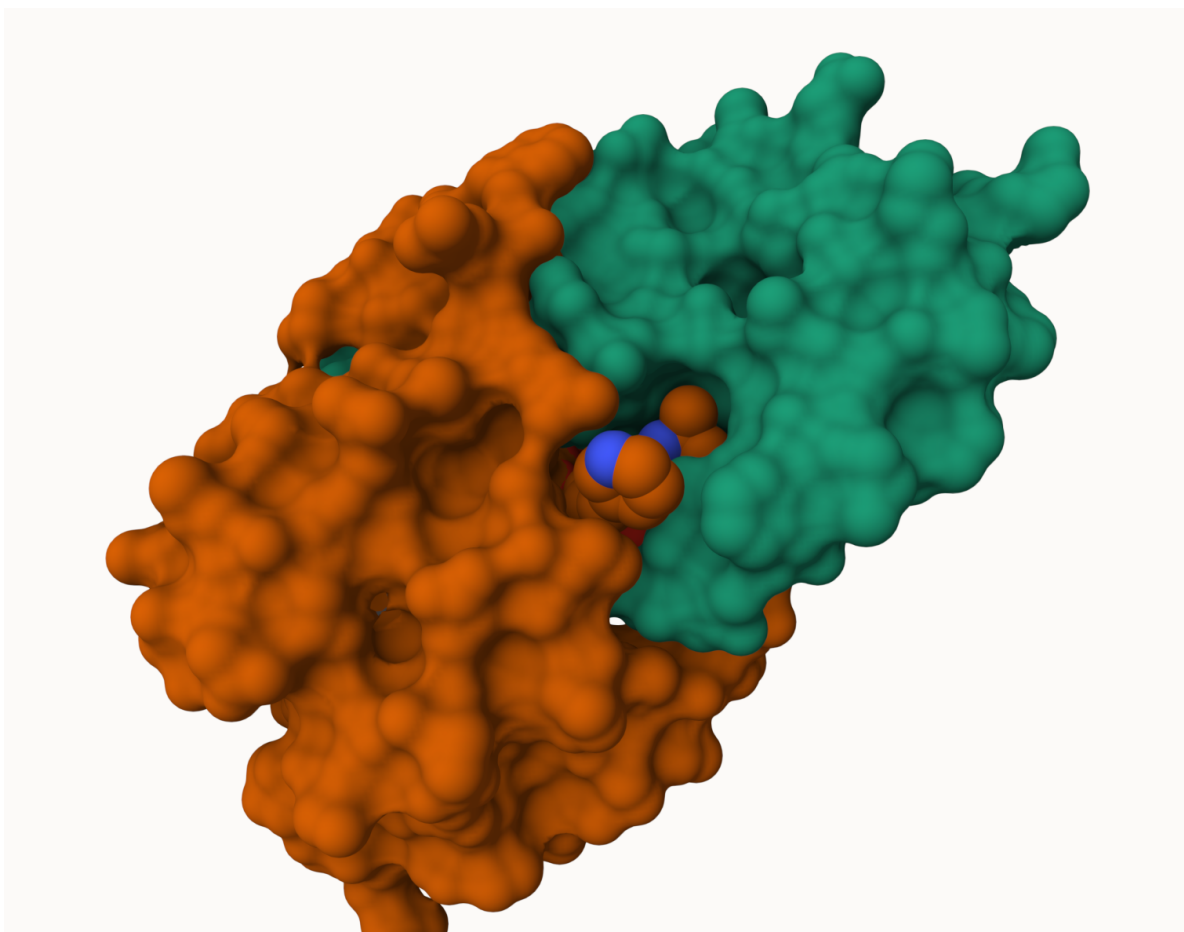


Figure 4: A third surface display image from molstar showing the Merk compound in the binding pocket

### The Bio3D package

The bio3d package allows us to do all sorts of structural bioinformatics work in R.

Let's start with how it can read these PDB files:

```
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
```

```
attributes(pdb)
```

```
$names
```

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

```
$class
```

```
[1] "pdb" "sse"
```

```
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>

```

```
pdbseq(pdb)
```

```

  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20
"P" "Q" "I" "T" "L" "W" "Q" "R" "P" "L" "V" "T" "I" "K" "I" "G" "G" "Q" "L" "K"
21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40
"E" "A" "L" "L" "D" "T" "G" "A" "D" "D" "T" "V" "L" "E" "E" "M" "S" "L" "P" "G"
41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60
"R" "W" "K" "P" "K" "M" "I" "G" "G" "I" "G" "G" "F" "I" "K" "V" "R" "Q" "Y" "D"
61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80
"Q" "I" "L" "I" "E" "I" "C" "G" "H" "K" "A" "I" "G" "T" "V" "L" "V" "G" "P" "T"
81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 1
"P" "V" "N" "I" "I" "G" "R" "N" "L" "L" "T" "Q" "I" "G" "C" "T" "L" "N" "F" "P"
  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20 21
"Q" "I" "T" "L" "W" "Q" "R" "P" "L" "V" "T" "I" "K" "I" "G" "G" "Q" "L" "K" "E"
22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41
"A" "L" "L" "D" "T" "G" "A" "D" "D" "T" "V" "L" "E" "E" "M" "S" "L" "P" "G" "R"
42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61
"W" "K" "P" "K" "M" "I" "G" "G" "I" "G" "G" "F" "I" "K" "V" "R" "Q" "Y" "D" "Q"
62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81
"I" "L" "I" "E" "I" "C" "G" "H" "K" "A" "I" "G" "T" "V" "L" "V" "G" "P" "T" "P"
82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99
"V" "N" "I" "I" "G" "R" "N" "L" "L" "T" "Q" "I" "G" "C" "T" "L" "N" "F"

```

Q7: How many amino acid residues are there in this pdb object?

```
sum(pdb$calpha)
```

```
[1] 198
```

There are 198 amino acid residues in this pdb object

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

The two non-protein resiudes are HOH and Mk1.

Q9: How many protein chains are in this structure?

There are two protein chains, chain A and B

```
unique(pdb$atom$chain)
```

```
[1] "A" "B"
```

##Predicting functional motions of a single structure

Let's do bioinformatics prediction of functional motion - i.e. the movements that one of these molecules needs to make to do its stuff.

```
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:

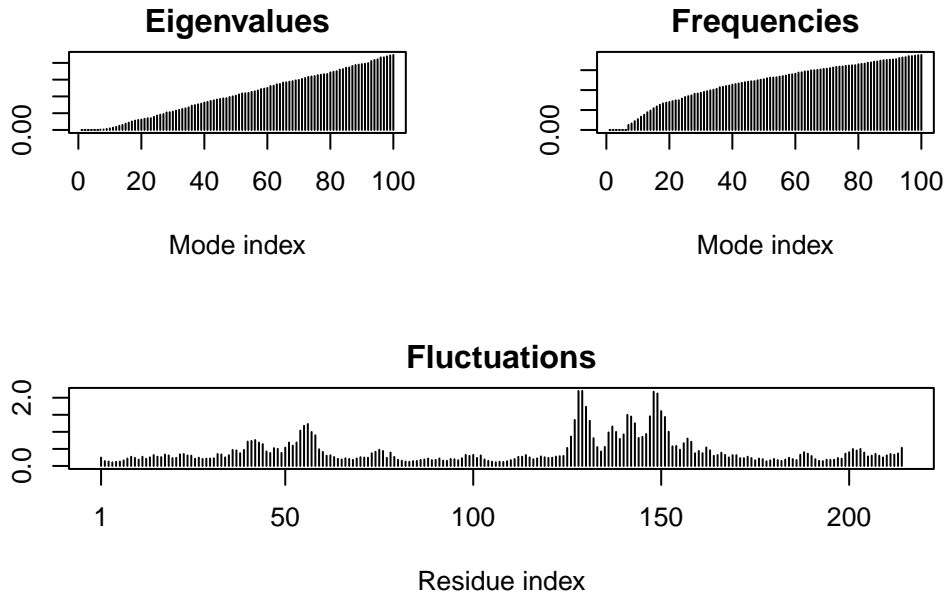
```
MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMRLRAAVKSGSELGKQAKDIMDAGKLV  
DELVIALVKERIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFDVPDELIVDKI  
VGRRVHAPSGRVYHVKFNPVKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQMTAPLIG  
YYSKEAEAGNTKYAKVDGTPVAEVRADLEKILG
```

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

```
# Perform flexibility prediction  
m <- nma(adk)
```

```
Building Hessian...      Done in 0.01 seconds.  
Diagonalizing Hessian... Done in 0.225 seconds.
```

```
plot(m)
```



We will write out multi-model trajectory PDB files that we can make animations of predicted motions from.

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

We can open this file in mol\* to play the trajectory.

## Comparative analysis of Protein structure

```
library(bio3d)
```

Here we will find and analyze all ADK structures in the PDB database.

We will start with the single database accession id: "1ake\_A"

```
id <- "1ake_A"
aa <- get.seq(id)
```

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

Fetching... Please wait. Done.

```
aa
```

```

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

      61      .      .      .      .      .      .      120
pdb|1AKE|A  DELVIALVKERIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFDVPDELIVDRI
      61      .      .      .      .      .      .      120

      121      .      .      .      .      .      .      180
pdb|1AKE|A  VGRRVHAPSGRVYHV KFNPPKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQM TAPLIG
      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

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

The `msa` package is found only on BioConductor.

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

The package “Grantlab/bio3d-view”

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

True, they are able to be used to install these packages.

Q13: How many amino acids are in this sequence, i.e. how long is this sequence?

```
ncol(aa$ali)
```

```
[1] 214
```

```
attributes(aa)
```

```
$names
```

```
[1] "id"    "ali"    "call"
```

```
$class
```

```
[1] "fasta"
```

There are 214 Amino acids.

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

```
#hits <- plot(b)
```

```
#hits
```

Pre-calculated results:

```
hits <- NULL
```

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

```
# Download releated PDB files
```

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

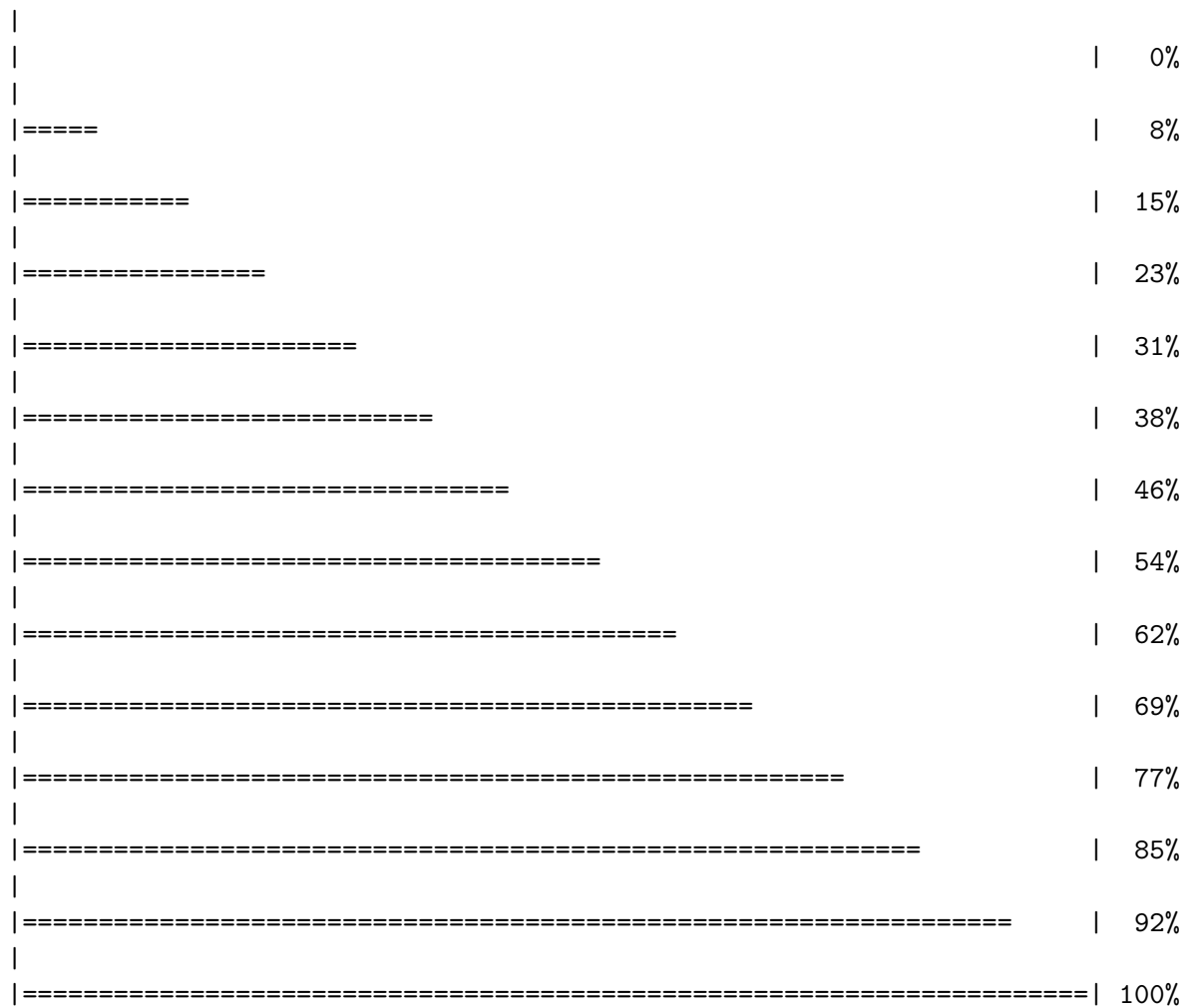
Warning in get.pdb(hits\$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):  
pdbs/6HAP.pdb.gz exists. Skipping download

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



Next we will use the `pdbaln()` function to align and also optionally fit the identified PDB structures.

```
# 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
```



```

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

```

```

[Truncated_Name:1] 1AKE_A.pdb
[Truncated_Name:2] 6S36_A.pdb
[Truncated_Name:3] 6RZE_A.pdb
[Truncated_Name:4] 3HPR_A.pdb
[Truncated_Name:5] 1E4V_A.pdb
[Truncated_Name:6] 5EJE_A.pdb
[Truncated_Name:7] 1E4Y_A.pdb
[Truncated_Name:8] 3X2S_A.pdb
[Truncated_Name:9] 6HAP_A.pdb
[Truncated_Name:10] 6HAM_A.pdb
[Truncated_Name:11] 4K46_A.pdb
[Truncated_Name:12] 3GMT_A.pdb
[Truncated_Name:13] 4PZL_A.pdb

1          .          .          .          40
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
-----MRIILLGAPVAGKGTQAQFIMEKYGIPQIS
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
-----MRIILLGALVAGKGTQAQFIMEKYGIPQIS
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
-----MRLILLGAPGAGKGTQANFIKEKFGIPQIS
-----TENLYFQSNAMRIILLGAPGAGKGTQAKIIEQKYNIAHIS
          **~*****  *****  *  *~ *  **
1          .          .          .          40

[Truncated_Name:1] 1AKE_A.pdb
[Truncated_Name:2] 6S36_A.pdb
[Truncated_Name:3] 6RZE_A.pdb
[Truncated_Name:4] 3HPR_A.pdb
[Truncated_Name:5] 1E4V_A.pdb
[Truncated_Name:6] 5EJE_A.pdb
[Truncated_Name:7] 1E4Y_A.pdb
[Truncated_Name:8] 3X2S_A.pdb
[Truncated_Name:9] 6HAP_A.pdb
[Truncated_Name:10] 6HAM_A.pdb
[Truncated_Name:11] 4K46_A.pdb
[Truncated_Name:12] 3GMT_A.pdb
[Truncated_Name:13] 4PZL_A.pdb

41          .          .          .          80
TGDMLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVKE
TGDMLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVKE
TGDMLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVKE
TGDMLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVKE
TGDMLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVKE
TGDMLRAAVKSGSELGKQAKDIMDACKLVDELVIALVKE
TGDMLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVKE
TGDMLRAAVKSGSELGKQAKDIMDCGKLVDELVIALVKE
TGDMLRAAVKSGSELGKQAKDIMDAGKLVDELVIALVRE
TGDMLRAAAIKSGSELGKQAKDIMDAGKLVDEIIIALVKE
TGDMLRAAAIKAGTELGKQAKSVIDAGQLVSDDIILGLVKE
TGDMLRAAVKAGTPLGVEAKTYMDEGKLVPSLIIGLVKE
TGDMIRETIKSGSALGQELKKVLDAGELVSDEFIIVKVD
****~*  ~* *~ **  *  ~*  ** *  ^^ ~* ^^
41          .          .          .          80

[Truncated_Name:1] 1AKE_A.pdb
[Truncated_Name:2] 6S36_A.pdb
[Truncated_Name:3] 6RZE_A.pdb
[Truncated_Name:4] 3HPR_A.pdb
[Truncated_Name:5] 1E4V_A.pdb

81          .          .          .          120
RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD

```

[Truncated_Name:6] 5EJE_A.pdb	RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:7] 1E4Y_A.pdb	RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:8] 3X2S_A.pdb	RIAQEDSRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:9] 6HAP_A.pdb	RICQEDSRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:10] 6HAM_A.pdb	RICQEDSRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:11] 4K46_A.pdb	RIAQDDCAKGFLLDGFPR TIPQADGLKEVGVVVDYVIEFD
[Truncated_Name:12] 3GMT_A.pdb	RLKEADCANGYLF DGFPR TIAQADAMKEAGVAIDYVLEID
[Truncated_Name:13] 4PZL_A.pdb	RISKNDCNNGFLLDGVPR TIPQAQELDKLGVNIDYIVEVD
	*^ * *~* ** ***** * ^ *~ ~**~* *
	81 . . . 120
	121 . . . 160
[Truncated_Name:1] 1AKE_A.pdb	VPDELIVDRIVGRRVHAPSGRVYHV KFNPPKVEGKDDVTG
[Truncated_Name:2] 6S36_A.pdb	VPDELIVDKIVGRRVHAPSGRVYHV KFNPPKVEGKDDVTG
[Truncated_Name:3] 6RZE_A.pdb	VPDELIVDAIVGRRVHAPSGRVYHV KFNPPKVEGKDDVTG
[Truncated_Name:4] 3HPR_A.pdb	VPDELIVDRIVGRRVHAPSGRVYHV KFNPPKVEGKDDGTG
[Truncated_Name:5] 1E4V_A.pdb	VPDELIVDRIVGRRVHAPSGRVYHV KFNPPKVEGKDDVTG
[Truncated_Name:6] 5EJE_A.pdb	VPDELIVDRIVGRRVHAPSGRVYHV KFNPPKVEGKDDVTG
[Truncated_Name:7] 1E4Y_A.pdb	VPDELIVDRIVGRRVHAPSGRVYHV KFNPPKVEGKDDVTG
[Truncated_Name:8] 3X2S_A.pdb	VPDELIVDRIVGRRVHAPSGRVYHV KFNPPKVEGKDDVTG
[Truncated_Name:9] 6HAP_A.pdb	VPDELIVDRIVGRRVHAPSGRVYHV KFNPPKVEGKDDVTG
[Truncated_Name:10] 6HAM_A.pdb	VPDELIVDRIVGRRVHAPSGRVYHV KFNPPKVEGKDDVTG
[Truncated_Name:11] 4K46_A.pdb	VADSVIVERMAGRRAHLASGR TYHNVPKVEGKDDVTG
[Truncated_Name:12] 3GMT_A.pdb	VPFSEIIERMSGRRTHPASGR TYHV KFNPPKVEGKDDVTG
[Truncated_Name:13] 4PZL_A.pdb	VADNLLIERITGRRIHPASGR TYHTKFNPPKVADKDDVTG
	* ^^^ ^ *** * *** * ~***** *** *
	121 . . . 160
	161 . . . 200
[Truncated_Name:1] 1AKE_A.pdb	EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:2] 6S36_A.pdb	EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:3] 6RZE_A.pdb	EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:4] 3HPR_A.pdb	EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:5] 1E4V_A.pdb	EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:6] 5EJE_A.pdb	EELTTRKDDQEECVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:7] 1E4Y_A.pdb	EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:8] 3X2S_A.pdb	EELTTRKDDQEETVRKRLCEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:9] 6HAP_A.pdb	EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:10] 6HAM_A.pdb	EELTTRKDDQEETVRKRLVEYHQMTAPLIGYYSKEAEAGN
[Truncated_Name:11] 4K46_A.pdb	EDLVIREDDKEETV LARLG VYHNQTAPLIAYYGKEAEAGN
[Truncated_Name:12] 3GMT_A.pdb	EPLVQRDDKEETVKKRLDVYEAQTKPLITYYGDWARRGA
[Truncated_Name:13] 4PZL_A.pdb	EPLITRTDDNEDTVKQRLSVYHAQTAKLIDFYRNFSSNT
	* * * * * ~* * * * * ~*

161 . . . 200

201 . . . 227

```
[Truncated_Name:1] 1AKE_A.pdb T--KYAKVDGTPVAEVRADLEKILG-
[Truncated_Name:2] 6S36_A.pdb T--KYAKVDGTPVAEVRADLEKILG-
[Truncated_Name:3] 6RZE_A.pdb T--KYAKVDGTPVAEVRADLEKILG-
[Truncated_Name:4] 3HPR_A.pdb T--KYAKVDGTPVAEVRADLEKILG-
[Truncated_Name:5] 1E4V_A.pdb T--KYAKVDGTPVAEVRADLEKILG-
[Truncated_Name:6] 5EJE_A.pdb T--KYAKVDGTPVAEVRADLEKILG-
[Truncated_Name:7] 1E4Y_A.pdb T--KYAKVDGTPVAEVRADLEKILG-
[Truncated_Name:8] 3X2S_A.pdb T--KYAKVDGTPVAEVRADLEKILG-
[Truncated_Name:9] 6HAP_A.pdb T--KYAKVDGTPVCEVRADLEKILG-
[Truncated_Name:10] 6HAM_A.pdb T--KYAKVDGTPVCEVRADLEKILG-
[Truncated_Name:11] 4K46_A.pdb T--QYLKFDGTKA VA EVSAE LEKALA-
[Truncated_Name:12] 3GMT_A.pdb E-----NGLKAPA-----YRKISG-
[Truncated_Name:13] 4PZL_A.pdb KIPKYIKINGDQAVEKVSQDIFDQLNK
```

\*

201 . . . 227

Call:

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

Class:

```
pdbs, fasta
```

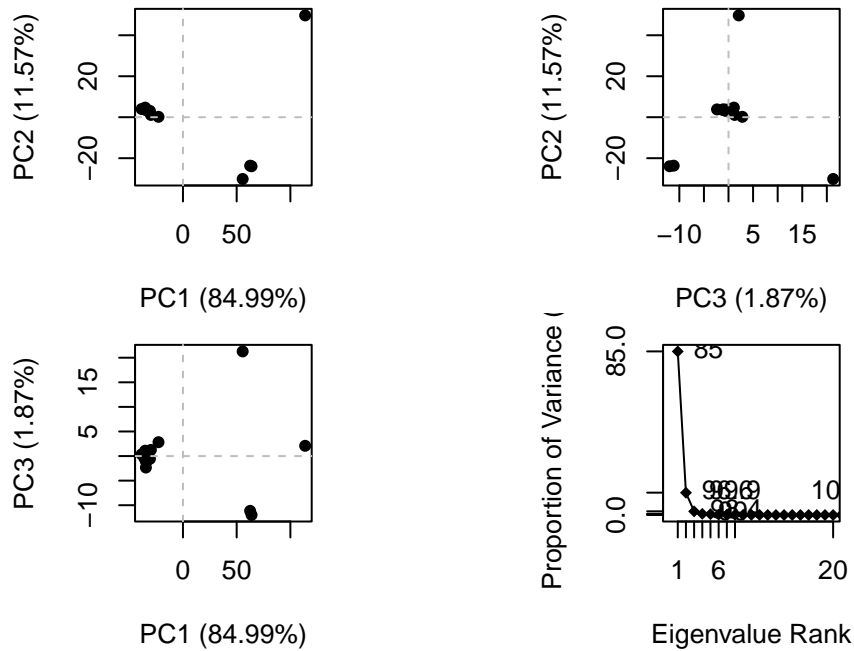
Alignment dimensions:

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

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

## Principal Component Analysis

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



To visualize the structural variations i will use the function `mktrj()` to generate a trajectory PDB file by interpolating along a give PC (eigenvector)

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

```
uniprot <- 24883887
pdb <- 195610
```

```
pdb/uniprot *100
```

```
[1] 0.0786091
```