

class10

james woolley A16440072

We should take a look into the PDB databse, which is the second oldest bioinformatics database in the world and the main repository of protein structures.

We should download the PDB data distribution from the website and move it to the folder.

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

	X.ray	EM	NMR	Multiple.methods	Neutron	Other
Protein (only)	158,844	11,759	12,296	197	73	32
Protein/Oligosaccharide	9,260	2,054	34	8	1	0
Protein/NA	8,307	3,667	284	7	0	0
Nucleic acid (only)	2,730	113	1,467	13	3	1
Other	164	9	32	0	0	0
Oligosaccharide (only)	11	0	6	1	0	4
Total						
Protein (only)	183,201					
Protein/Oligosaccharide	11,357					
Protein/NA	12,265					
Nucleic acid (only)	4,327					
Other	205					
Oligosaccharide (only)	22					

Currently, these values are stored as cahracters because there are commas separating the thousands place. We can fix this using the gsub() function.

```
x <- stats$X.ray #this will use the gsub function to globally substitute commas with nothing
as.numeric(gsub(",", "", x))
```

```
[1] 158844 9260 8307 2730 164 11
```

```
rm.comma <- function(x){#we can write a function to remove commas
  as.numeric( gsub(",", "", x))
}
```

```
rm.comma(stats$EM)
```

```
[1] 11759 2054 3667 113 9 0
```

```
pdbstats <- apply(stats, 2, rm.comma)
```

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

We can do some math to find out how many structures are solved by X-ray or EM

```
totals <- apply(pdbstats, 2, sum)
round(totals/totals["Total"] *100, 2)
```

X-ray	EM	NMR	Multiple.methods
84.83	8.33	6.68	0.11
Neutron	Other	Total	
0.04	0.02	100.00	

We can see that 84.83% of the structures were determined by X-ray and 8.33% of the structures were determined by EM. >Q2 What proportion of structures in the PDB are protein?

```
round(pdbstats[1,"Total"]/sum(pdbstats), 2)
```

```
Total
0.43
```

A total of 43% of the structures are protein. >Q3.

We skipped this one <3

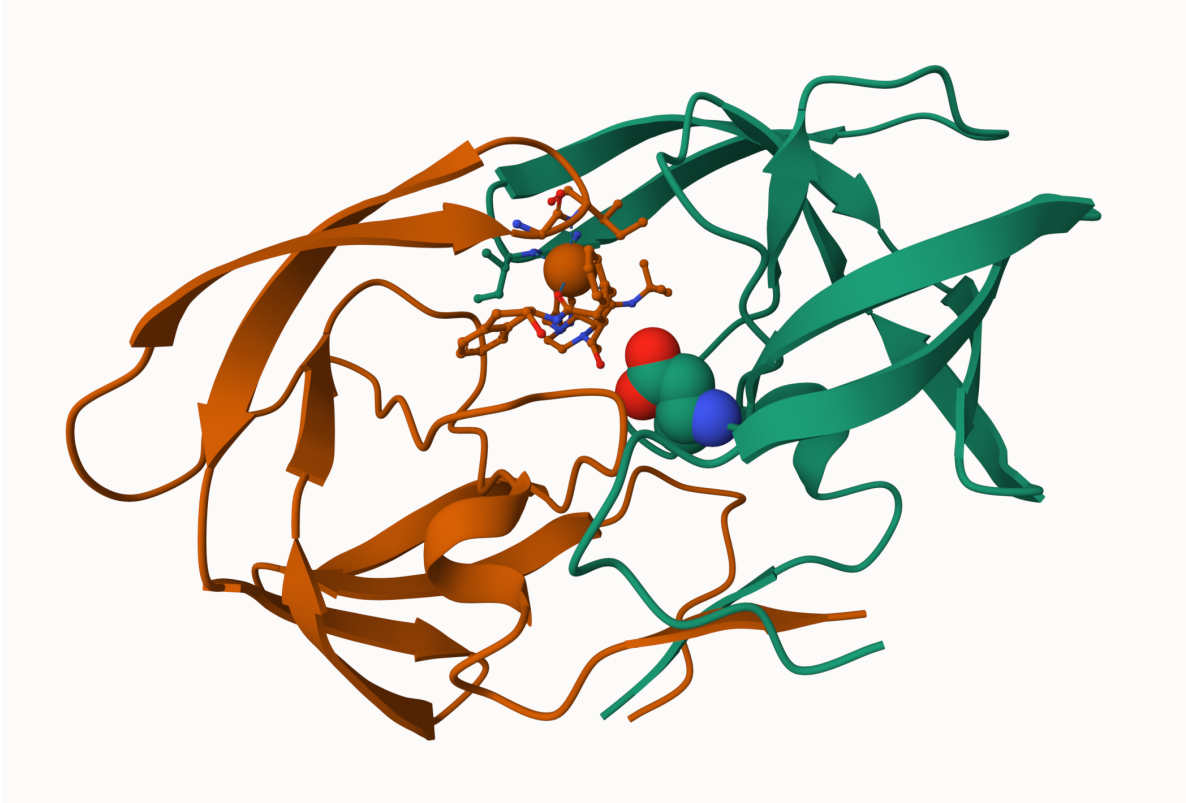
Q4. Water molecules normally have 3 atoms. Why do we see just one atom per water molecule in this structure?

The hydrogens don't appear because they're smaller than the resolution the protein was analyzed with (2 amstrongs)

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?

H308 is the critical water molecule that allows the protein to work. It has been identified in the figure below.

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.



```
##The bio3d package for structural bioinformatics
```

```
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
```

```
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>										

#predicting functional motions

We can finish today with bioinformatics calculations to predict functional motions of a PDB 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:

```
MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMLRAAVKSGSELGKQAKDIMDAGKLV  
TDELVIALVKERIAQEDCRNGFLDGFPR TIPQADAMKEAGINVDYVLEFDVPDELIVDKI  
VGRRVHAPSGRVYHVKFNPPKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQM TAPLIG  
YYSKEAEAGNTKYAKVDGTPVAEVRADLEKILG
```

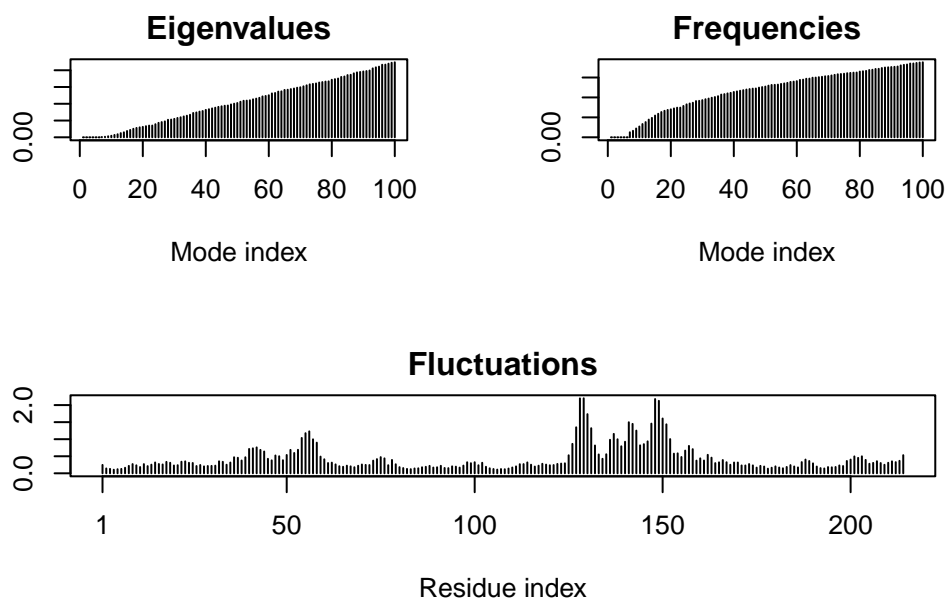
+ attr: atom, xyz, seqres, helix, sheet,
calpha, remark, call

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

Building Hessian... Done in 0.014 seconds.

Diagonalizing Hessian... Done in 0.271 seconds.

```
plot(m)
```



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

We can use the pdb file we generated and put it into molstar to visualise how it moves.