# lab9

Meha Thakur PID 16020450

## Table of contents

Reading in files and loading packages

```
library(tidyverse)
```

```
-- Attaching core tidyverse packages ---------------------- tidyverse 2.0.0 --
v dplyr     1.1.4     v readr     2.1.5
v forcats   1.0.1     v stringr   1.6.0
v ggplot2   4.0.0     v tibble    3.3.0
v lubridate 1.9.4     v tidyr     1.3.1
v purrr     1.0.2
-- Conflicts ---------------------------------------- tidyverse_conflicts() --
x dplyr::filter() masks stats::filter()
x dplyr::lag()    masks stats::lag()
i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to becom
```

```
library(readr)

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

```
#one way to change characters into numerics - annoying way
num_data<- as.numeric(sub(",","",data$X.ray))

#other way using readr to import data
```

```r
data1<-read_csv("Data Export Summary.csv")
```

```
Rows: 6 Columns: 9
-- Column specification ---------------------------------------------------
Delimiter: ","
chr (1): Molecular Type
dbl (4): Integrative, 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.
```

```r
XrayT<-sum(data1$"X-ray")
emT<- sum(data1$EM)
total<-sum(data1$Total)

round(XrayT/total *100,2) #x ray proportion
```

```
[1] 81.43
```

```r
round(emT/total*100,2) #EM proportion
```

```
[1] 12.27
```

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

ans: 81.43% X-ray, 12.27% EM

Q2: What proportion of structures in the PDB are protein?
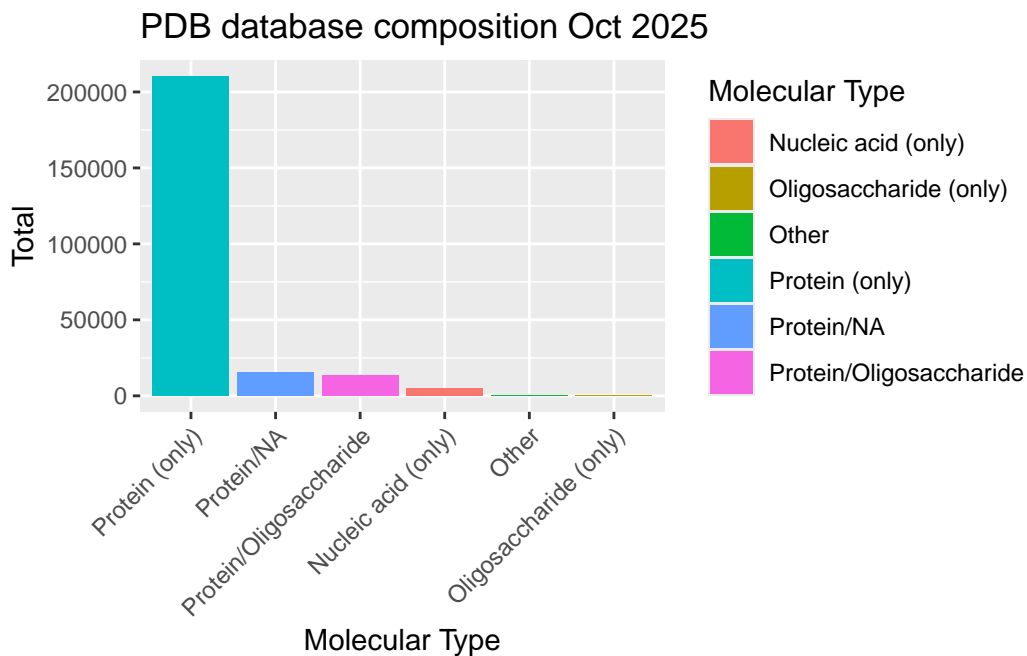
```r
prot.total<-sum(data1$Total[1])

round(prot.total/total*100,2)
```

```
[1] 86.05
```

ans: 86.05%

plot the PDB database composition oct 2025. Want stacked barchart by number in each column. Want data in long format,

```
#regular bar chart
ggplot(data1,aes(x=reorder(`Molecular Type`,-Total),y=Total,fill=`Molecular Type`))+
  geom_col()+
  labs(x="Molecular Type",y="Total", title = "PDB database composition Oct 2025")+
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
```



```
#stacked bar chart -- NOT DONE

#need to convert data into long format
data_long <- data1 %>%
  select(-Total)%>%
  pivot_longer(
    cols = where(is.numeric),
    names_to = "Experimental Method",
    values_to = "tot.val"
  )

ggplot(data_long,aes(x=reorder(`Molecular Type`,-tot.val),y=tot.val,fill=`Experimental Method
```
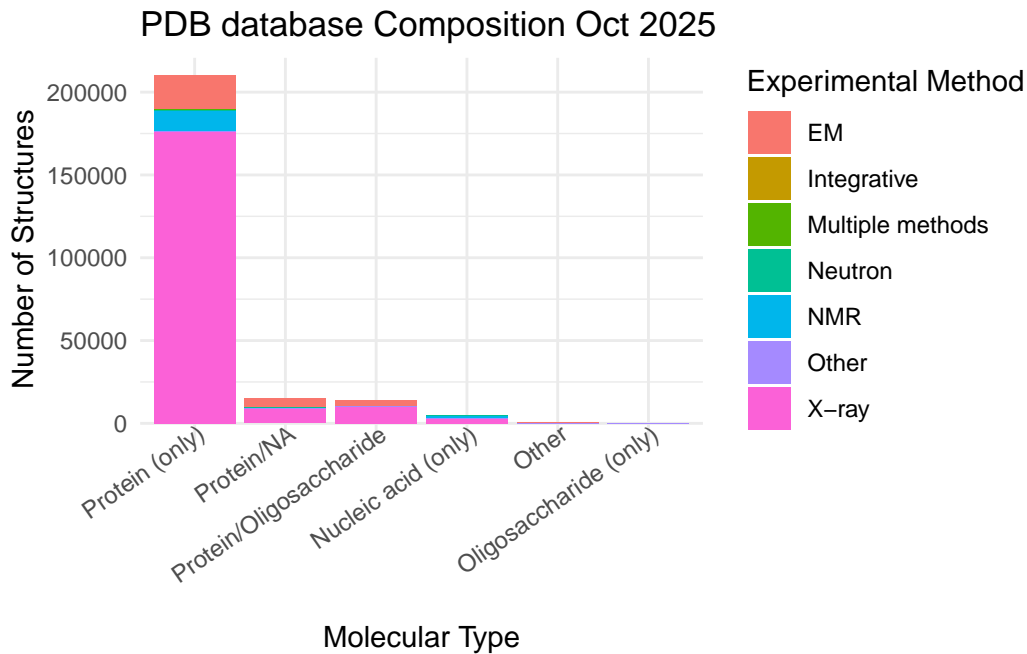
```
geom_col(position="stack")+
labs(x="Molecular Type",y="Number of Structures", title = "PDB database Composition Oct 202
theme_minimal() +
theme(axis.text.x = element_text(angle = 35,hjust=1,vjust=1.2))
```

## PDB database Composition Oct 2025



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?

HIV search came up with 4886 structures.

### Visualising structure data

The Mol* Viewer is embeded in many bioinfo websites. Homepage is https://molstar.org/

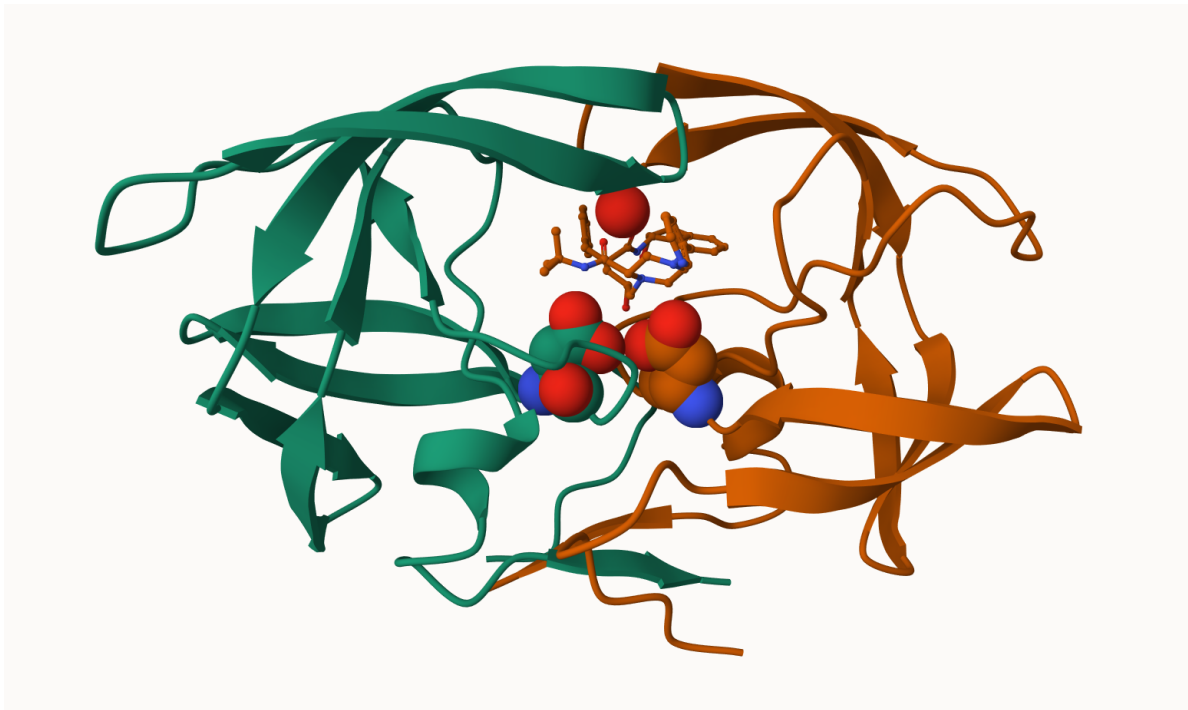We can insert any figure/image file using markdown This is HIV Protease

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

It would look too crowded to have each individual atom in the water molecule on this 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

Water 308, highlighted above as a red circle in the image below near the highlighted amino acids.

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. Discussion Topic: Can you think of a way in which indinavir, or even larger ligands and substrates, could enter the binding site?



## Bio3D package for structural informatics

We can use the bio3d package to read and analuse biomolecular data in R:

```r
library(bio3d)

#read in pdb file
hiv<-read.pdb("1hsg")
```

  Note: Accessing on-line PDB file

```r
#preview atom
head(hiv$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>
```

```r
#sequence - chain A and B
pdbseq(hiv)
```

```
  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"
```

```
#trim to chain A
chainA<-trim.pdb(hiv, chain='A')
chainA.seq<-pdbseq(chainA)

hiv
```

```
 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?

198 residues

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

HOH, MK1

Q9: How many protein chains are in this structure?

2 chains

Blast

```
#blast<- blast.pdb(chainA.seq) commented out as I'm doing it below in cached results chunk
```

```
#when I render this, it will use the saved value and not have to rerun the blast code again
blast<- blast.pdb(chainA.seq)
```
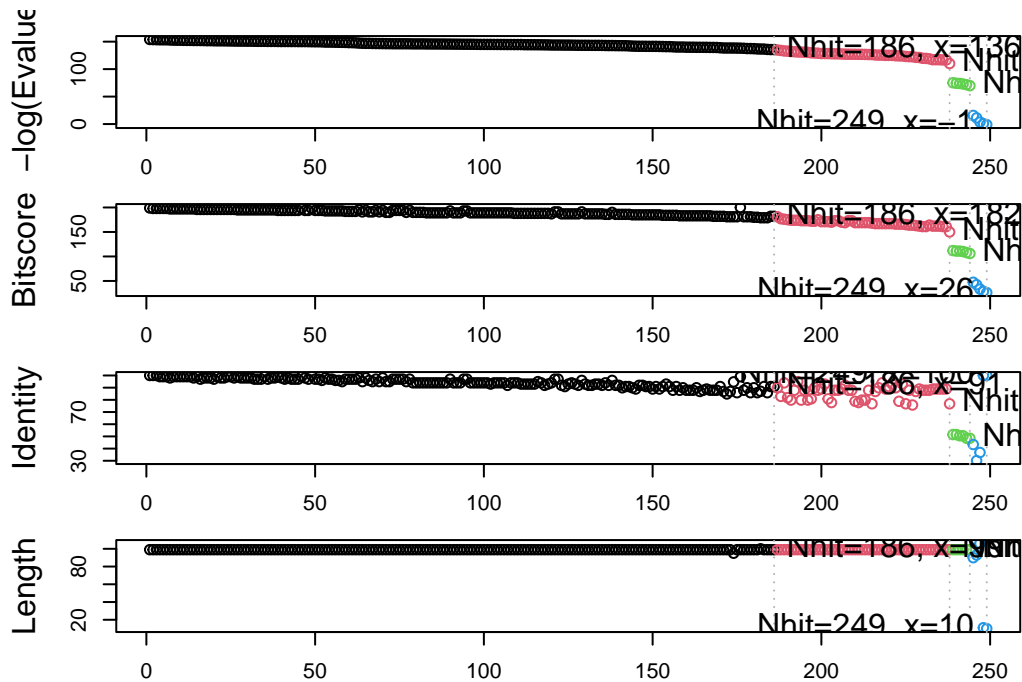
```
 Searching ... please wait (updates every 5 seconds) RID = G56ZGT3N014

 Reporting 249 hits
```

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

```
 * Possible cutoff values:    135 110 69 -2
           Yielding Nhits:    186 238 244 249

 * Chosen cutoff value of:    69
           Yielding Nhits:    244
```

```
hits$pdb.id #gets all the accession numbers of those 249 hits
```

```
  [1] "1W5V_A" "2FDE_A" "1AJV_A" "2R38_A" "2R3T_A" "1HXB_A" "1BV9_A" "1AAQ_A"
  [9] "1AXA_A" "1HVS_A" "1ZP8_A" "2QHC_A" "1A8G_A" "2O4L_A" "5COK_A" "1TCX_A"
 [17] "2Z54_A" "1D4S_A" "1BV7_A" "1BWA_A" "1A9M_A" "2FLE_A" "1ODY_A" "1GNN_A"
 [25] "1GNM_A" "5YRS_B" "1HEF_E" "1ODX_A" "4QGI_A" "1BVE_A" "2AZ8_A" "1A30_A"
 [33] "6DH6_A" "6DH0_A" "2I4D_A" "6OOS_A" "1RL8_A" "5YRS_A" "1ZSF_A" "2Q64_A"
 [41] "6DH3_A" "2NPH_A" "2Q63_A" "1LZQ_A" "1FB7_A" "1G6L_A" "1HIV_A" "6OOU_A"
 [49] "1HVC_A" "2I4V_A" "2AZ9_A" "6OOT_A" "2P3B_B" "5KAO_A" "2WL0_A" "6OPT_A"
 [57] "1IZI_A" "1MRX_A" "2PYM_A" "2PYN_A" "1DMP_A" "4K4P_A" "1LV1_A" "1AID_A"
 [65] "1LV1_A" "1ZBG_A" "3TKG_A" "1HVC_A" "5YOK_A" "1G6L_A" "1FGC_C" "3K4V_A"
 [73] "3KT5_A" "3KT5_A" "4QLH_A" "4QLH_A" "2F3K_A" "4Q5M_A" "2AOC_A" "3B80_A"
 [81] "3VF5_A" "2AVQ_A" "1DW6_C" "1KZK_A" "2HS1_A" "1K6C_A" "1MTB_A" "4Q1X_A"
 [89] "4Q1W_A" "4Q5M_A" "3D1X_A" "2AVM_A" "3PWM_A" "3KT2_A" "3KT2_A" "1SDV_A"
 [97] "3JVW_A" "3OY4_A" "1A94_A" "2HS2_A" "4EJ8_A" "2FGU_A" "2AVV_A" "3JW2_A"
[105] "3BVA_A" "1FFF_C" "3S43_B" "2NXD_A" "1FG6_C" "1EBK_C" "4Q1Y_A" "3EL4_A"
[113] "1F7A_A" "1K2B_A" "2FGV_A" "1Z8C_A" "2G69_A" "3EL9_A" "3OXV_A" "1BDR_A"
[121] "3N3I_A" "3N3I_A" "3OXW_A" "3S43_A" "3EM3_A" "3CYW_A" "5KQX_A" "2B60_A"
[129] "7DOZ_A" "1K2C_A" "1MT7_A" "3EM4_A" "4QJ9_A" "1BDL_A" "3LZS_A" "5T84_A"
[137] "4DQB_A" "7DOZ_A" "4QJ2_A" "3LZV_A" "1SGU_A" "2FXE_A" "1BDQ_A" "3U71_A"
[145] "2R5P_A" "4OBD_A" "7MAS_A" "3IXO_A" "3D3T_A" "5YOJ_A" "3LZU_A" "4NJS_A"
[153] "3EKP_A" "1B6J_A" "3EKQ_A" "2RKF_A" "1C6X_A" "7MAR_A" "4DQF_A" "1RPI_A"
[161] "3OU1_B" "3PJ6_A" "2P3A_A" "6OGQ_A" "3OQ7_A" "5KR1_A" "3OQD_A" "4RVI_A"
[169] "3OQA_A" "1B6K_A" "3OUD_B" "6MK9_A" "3S09_A" "1Q9P_A" "6I45_A" "7SEP_A"
[177] "4NJT_A" "3BXR_A" "4YOA_A" "4DQC_A" "2FDD_A" "2RKG_A" "4DQH_A" "2P3C_A"
[185] "4EP2_A" "4EP2_A" "4EQ0_A" "4NPT_A" "6OPU_A" "4NPU_A" "3U7S_A" "3HAW_A"
[193] "2AZB_A" "3TTP_A" "3HB0_A" "3GGU_A" "7N6T_A" "6OPV_A" "4EQ0_A" "6OPX_A"
[201] "2O4N_A" "5T2E_A" "3UCB_A" "3KA2_A" "3FSM_A" "6OPW_A" "2AZC_A" "3FSM_A"
[209] "3HLO_A" "2P3D_A" "3T3C_A" "7MYP_A" "6O54_X" "6OPY_A" "4Z4X_A" "6OPZ_A"
[217] "2JE4_A" "1DAZ_C" "7MAP_A" "7MAQ_A" "1K1U_A" "2B7Z_A" "3MWS_A" "1K1T_A"
[225] "8DCH_A" "3I2L_A" "6P9A_A" "2FXD_A" "2J9J_A" "3DCK_A" "2J9J_B" "3NXE_A"
[233] "2O40_A" "2O40_A" "3NXE_A" "3KA2_A" "3HLO_A" "5B18_A" "1SIP_A" "2SAM_A"
[241] "1AZ5_A" "1SIV_A" "1HII_A" "1IVP_A"
```

## Prediction of Functional Motions

We can run Normal Mode Analysis (NMA) to predict large scale motions/flexibility/dynamics of any biomolecule that we can read into R.

Lets look into ADK

```
adk<-read.pdb("1ake")
```

```
 Note: Accessing on-line PDB file
   PDB has ALT records, taking A only, rm.alt=TRUE
```

```
adk_A<-trim.pdb(adk, chain="A")
adk_A
```

```
 Call:  trim.pdb(pdb = adk, chain = "A")

   Total Models#: 1
     Total Atoms#: 1954,  XYZs#: 5862  Chains#: 1  (values: A)

     Protein Atoms#: 1656  (residues/Calpha atoms#: 214)
     Nucleic acid Atoms#: 0  (residues/phosphate atoms#: 0)

     Non-protein/nucleic Atoms#: 298  (residues: 242)
     Non-protein/nucleic resid values: [ AP5 (1), HOH (241) ]

   Protein sequence:
     MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMLRAAVKSGSELGKQAKDIMDAGKLVT
     DELVIALVKERIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFDVPDELIVDRI
     VGRRVHAPSGRVYHVKFNPPKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQMTAPLIG
     YYSKEAEAGNTKYAKVDGTKPVAEVRADLEKILG

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

```
m<-nma(adk_A)
```

```
 Building Hessian...        Done in 0.084 seconds.
 Diagonalizing Hessian...   Done in 0.249 seconds.
```

```
plot(m)
```

Lets write out a "trajectory" of predicted motion

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

## Play with 3D viewing in R

We can use the new **bio3dview** package, which is not yet on CRAN, to render interactive 3D views in R and HTML quarto output reports. Need to install from GitHub. We can use the **pak** package. I am choosing to use **devtools** as is outlined on the lab walk through.

```
library(BiocManager)
```

```
Bioconductor version '3.16' is out-of-date; the current release version '3.22'
  is available with R version '4.5'; see https://bioconductor.org/install
```

```
library(devtools)
```

```
Loading required package: usethis
```

```
Attaching package: 'devtools'
```

The following object is masked from 'package:BiocManager':

    install

```r
library(bio3d)
library(bio3d.view)
```

Attaching package: 'bio3d.view'

The following object is masked from 'package:tibble':

    view

```r
library(msa)
```

Loading required package: Biostrings

Loading required package: BiocGenerics

Attaching package: 'BiocGenerics'

The following objects are masked from 'package:lubridate':

    intersect, setdiff, union

The following objects are masked from 'package:dplyr':

    combine, intersect, setdiff, union

The following objects are masked from 'package:stats':

    IQR, mad, sd, var, xtabs

```
The following objects are masked from 'package:base':

    anyDuplicated, aperm, append, as.data.frame, basename, cbind,
    colnames, dirname, do.call, duplicated, eval, evalq, Filter, Find,
    get, grep, grepl, intersect, is.unsorted, lapply, Map, mapply,
    match, mget, order, paste, pmax, pmax.int, pmin, pmin.int,
    Position, rank, rbind, Reduce, rownames, sapply, setdiff, sort,
    table, tapply, union, unique, unsplit, which.max, which.min


Loading required package: S4Vectors


Loading required package: stats4



Attaching package: 'S4Vectors'


The following objects are masked from 'package:lubridate':

    second, second<-


The following objects are masked from 'package:dplyr':

    first, rename


The following object is masked from 'package:tidyr':

    expand


The following objects are masked from 'package:base':

    expand.grid, I, unname


Loading required package: IRanges



Attaching package: 'IRanges'


The following object is masked from 'package:bio3d':

    trim
```

```
The following object is masked from 'package:lubridate':

    %within%

The following objects are masked from 'package:dplyr':

    collapse, desc, slice

The following object is masked from 'package:purrr':

    reduce

Loading required package: XVector


Attaching package: 'XVector'

The following object is masked from 'package:purrr':

    compact

Loading required package: GenomeInfoDb


Attaching package: 'Biostrings'

The following object is masked from 'package:bio3d':

    mask

The following object is masked from 'package:base':

    strsplit


Attaching package: 'msa'

The following object is masked from 'package:BiocManager':

    version
```

```
#library("Grantlab/bio3d-view")
```

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

msa, as it installed using biocmanager

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

bio3d.view

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

TRUE

```
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   MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMLRAAVKSGSELGKQAKDIMDAGKLVT
            1        .         .         .         .         .        60

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

            121      .         .         .         .         .        180
pdb|1AKE|A   VGRRVHAPSGRVYHVKFNPPKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQMTAPLIG
            121      .         .         .         .         .        180

            181      .         .          .    214
pdb|1AKE|A   YYSKEAEAGNTKYAKVDGTKPVAEVRADLEKILG
            181      .         .          .    214

Call:
  read.fasta(file = outfile)
```

16

```
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

Commented out as document won't render.

```
# Blast or hmmer search
#b <- blast.pdb(aa)
```

also commented out as above is not working

```
# Plot a summary of search results
#hits <- plot(b)
# List out some 'top hits'
#head(hits$pdb.id)
```

putting these into a vector

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

This chunk works, however will not render unless I comment out.

```
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%
```

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

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

Sequence Alignment Overview

annotating

```
anno <- pdb.annotate(ids)
unique(anno$source)
```

```
[1] "Escherichia coli"
[2] "Escherichia coli K-12"
[3] "Escherichia coli O139:H28 str. E24377A"
[4] "Escherichia coli str. K-12 substr. MDS42"
[5] "Photobacterium profundum"
[6] "Burkholderia pseudomallei 1710b"
[7] "Francisella tularensis subsp. tularensis SCHU S4"
```

```
        structureId chainId macromoleculeType chainLength experimentalTechnique
1AKE_A         1AKE       A          Protein         214                X-ray
6S36_A         6S36       A          Protein         214                X-ray
6RZE_A         6RZE       A          Protein         214                X-ray
3HPR_A         3HPR       A          Protein         214                X-ray
1E4V_A         1E4V       A          Protein         214                X-ray
5EJE_A         5EJE       A          Protein         214                X-ray
1E4Y_A         1E4Y       A          Protein         214                X-ray
3X2S_A         3X2S       A          Protein         214                X-ray
6HAP_A         6HAP       A          Protein         214                X-ray
6HAM_A         6HAM       A          Protein         214                X-ray
4K46_A         4K46       A          Protein         214                X-ray
3GMT_A         3GMT       A          Protein         230                X-ray
4PZL_A         4PZL       A          Protein         242                X-ray
        resolution        scopDomain                                       pfam
1AKE_A        2.00 Adenylate kinase Adenylate kinase, active site lid (ADK_lid)
6S36_A        1.60              <NA> Adenylate kinase, active site lid (ADK_lid)
6RZE_A        1.69              <NA> Adenylate kinase, active site lid (ADK_lid)
3HPR_A        2.00              <NA>                     Adenylate kinase (ADK)
1E4V_A        1.85 Adenylate kinase Adenylate kinase, active site lid (ADK_lid)
5EJE_A        1.90              <NA>                     Adenylate kinase (ADK)
1E4Y_A        1.85 Adenylate kinase                     Adenylate kinase (ADK)
3X2S_A        2.80              <NA>                     Adenylate kinase (ADK)
6HAP_A        2.70              <NA> Adenylate kinase, active site lid (ADK_lid)
6HAM_A        2.55              <NA>                     Adenylate kinase (ADK)
4K46_A        2.01              <NA> Adenylate kinase, active site lid (ADK_lid)
3GMT_A        2.10              <NA> Adenylate kinase, active site lid (ADK_lid)
4PZL_A        2.10              <NA> Adenylate kinase, active site lid (ADK_lid)
               ligandId
1AKE_A             AP5
6S36_A CL (3),NA,MG (2)
6RZE_A     NA (3),CL (2)
3HPR_A             AP5
1E4V_A             AP5
5EJE_A          AP5,CO
1E4Y_A             AP5
3X2S_A   JPY (2),AP5,MG
6HAP_A             AP5
6HAM_A             AP5
4K46_A     ADP,AMP,PO4
```

```
3GMT_A            SO4 (2)
4PZL_A        CA,FMT,GOL
                                                                      ligandName
1AKE_A                                    BIS(ADENOSINE)-5'-PENTAPHOSPHATE
6S36_A                   CHLORIDE ION (3),SODIUM ION,MAGNESIUM ION (2)
6RZE_A                             SODIUM ION (3),CHLORIDE ION (2)
3HPR_A                                    BIS(ADENOSINE)-5'-PENTAPHOSPHATE
1E4V_A                                    BIS(ADENOSINE)-5'-PENTAPHOSPHATE
5EJE_A               BIS(ADENOSINE)-5'-PENTAPHOSPHATE,COBALT (II) ION
1E4Y_A                                    BIS(ADENOSINE)-5'-PENTAPHOSPHATE
3X2S_A N-(pyren-1-ylmethyl)acetamide (2),BIS(ADENOSINE)-5'-PENTAPHOSPHATE,MAGNESIUM ION
6HAP_A                                    BIS(ADENOSINE)-5'-PENTAPHOSPHATE
6HAM_A                                    BIS(ADENOSINE)-5'-PENTAPHOSPHATE
4K46_A          ADENOSINE-5'-DIPHOSPHATE,ADENOSINE MONOPHOSPHATE,PHOSPHATE ION
3GMT_A                                                      SULFATE ION (2)
4PZL_A                                    CALCIUM ION,FORMIC ACID,GLYCEROL
                                         source
1AKE_A                        Escherichia coli
6S36_A                        Escherichia coli
6RZE_A                        Escherichia coli
3HPR_A                   Escherichia coli K-12
1E4V_A                        Escherichia coli
5EJE_A         Escherichia coli O139:H28 str. E24377A
1E4Y_A                        Escherichia coli
3X2S_A       Escherichia coli str. K-12 substr. MDS42
6HAP_A         Escherichia coli O139:H28 str. E24377A
6HAM_A                   Escherichia coli K-12
4K46_A                  Photobacterium profundum
3GMT_A              Burkholderia pseudomallei 1710b
4PZL_A Francisella tularensis subsp. tularensis SCHU S4

1AKE_A STRUCTURE OF THE COMPLEX BETWEEN ADENYLATE KINASE FROM ESCHERICHIA COLI AND THE INHIB:
6S36_A
6RZE_A
3HPR_A
1E4V_A
5EJE_A                                                                              Cryst
1E4Y_A
3X2S_A
6HAP_A
6HAM_A
4K46_A
3GMT_A
```

```
4PZL_A                                                                    The cryst
                                        citation rObserved    rFree
1AKE_A            Muller, C.W., et al. J Mol Biology (1992)    0.19600      NA
6S36_A               Rogne, P., et al. Biochemistry (2019)    0.16320 0.23560
6RZE_A               Rogne, P., et al. Biochemistry (2019)    0.18650 0.23500
3HPR_A  Schrank, T.P., et al. Proc Natl Acad Sci U S A (2009) 0.21000 0.24320
1E4V_A              Muller, C.W., et al. Proteins (1993)      0.19600      NA
5EJE_A  Kovermann, M., et al. Proc Natl Acad Sci U S A (2017) 0.18890 0.23580
1E4Y_A              Muller, C.W., et al. Proteins (1993)      0.17800      NA
3X2S_A             Fujii, A., et al. Bioconjug Chem (2015)    0.20700 0.25600
6HAP_A           Kantaev, R., et al. J Phys Chem B (2018)     0.22630 0.27760
6HAM_A           Kantaev, R., et al. J Phys Chem B (2018)     0.20511 0.24325
4K46_A              Cho, Y.-J., et al. To be published        0.17000 0.22290
3GMT_A Buchko, G.W., et al. Biochem Biophys Res Commun (2010) 0.23800 0.29500
4PZL_A                  Tan, K., et al. To be published       0.19360 0.23680
          rWork spaceGroup
1AKE_A 0.19600  P 21 2 21
6S36_A 0.15940    C 1 2 1
6RZE_A 0.18190    C 1 2 1
3HPR_A 0.20620  P 21 21 2
1E4V_A 0.19600  P 21 2 21
5EJE_A 0.18630  P 21 2 21
1E4Y_A 0.17800   P 1 21 1
3X2S_A 0.20700 P 21 21 21
6HAP_A 0.22370    I 2 2 2
6HAM_A 0.20311       P 43
4K46_A 0.16730 P 21 21 21
3GMT_A 0.23500   P 1 21 1
4PZL_A 0.19130       P 32
```
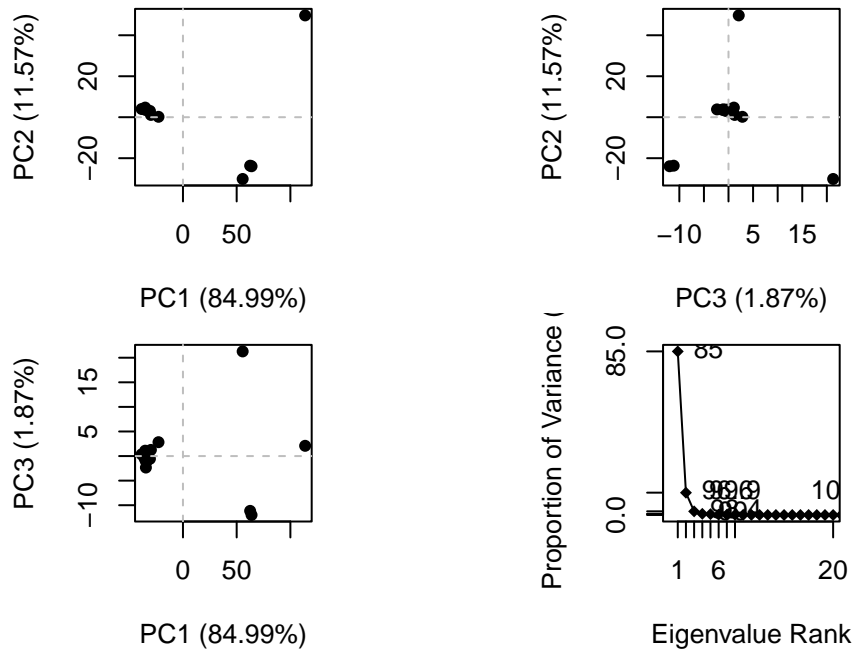
PCA

```
# Perform PCA
pc.xray <- pca(pdbs)
plot(pc.xray)
```
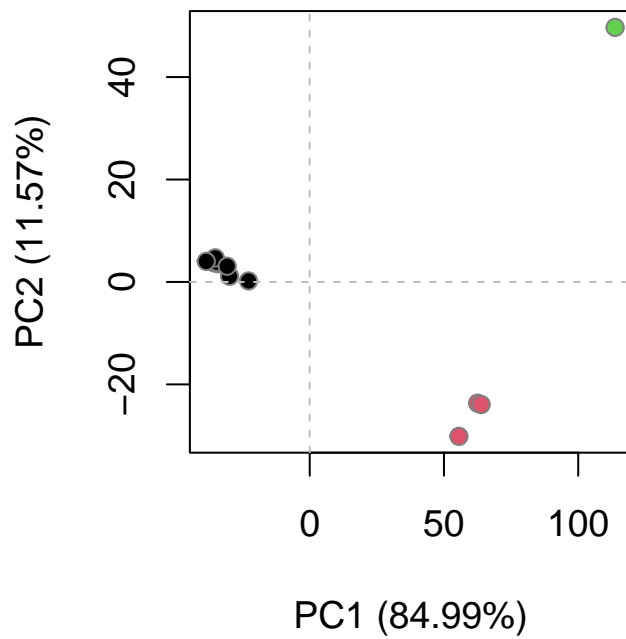
```
# Calculate RMSD
rd <- rmsd(pdbs)
```

Warning in rmsd(pdbs): 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)
```

## 11/05 addition: Comparitive analysis

```
library (bio3d)

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    MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMLRAAVKSGSELGKQAKDIMDAGKLVT
              1        .         .         .         .         .        60

              61       .         .         .         .         .       120
pdb|1AKE|A    DELVIALVKERIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFDVPDELIVDRI
```

```
           61         .         .         .         .         .        120


          121         .         .         .         .         .        180
pdb|1AKE|A     VGRRVHAPSGRVYHVKFNPPKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQMTAPLIG
          121         .         .         .         .         .        180


          181         .         .         .    214
pdb|1AKE|A     YYSKEAEAGNTKYAKVDGTKPVAEVRADLEKILG
          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
```

```r
#blast<- blast.pdb(aa)
#hits<- plot(blast)
hits <- NULL
hits$pdb.id <- c('1AKE_A','6S36_A','6RZE_A','3HPR_A','1E4V_A','5EJE_A','1E4Y_A','3X2S_A','6H
```

```r
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"
```

Download all these "hits" that are similar to our starting ID sequence

```r
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%
  |
```
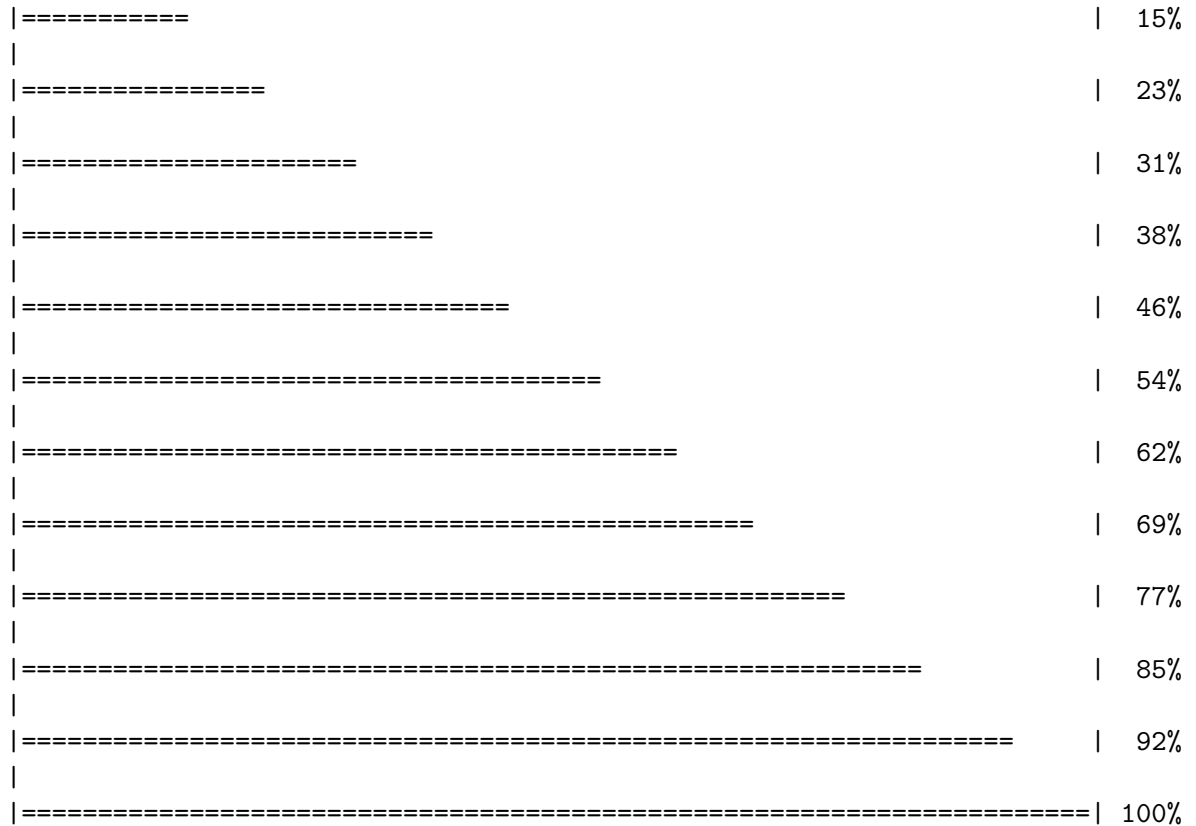
```
|==========                                                         |  15%
|
|===============                                                    |  23%
|
|====================                                               |  31%
|
|=========================                                          |  38%
|
|==============================                                     |  46%
|
|===================================                                |  54%
|
|========================================                           |  62%
|
|=============================================                      |  69%
|
|==================================================                 |  77%
|
|=======================================================            |  85%
|
|============================================================       |  92%
|
|==================================================================| 100%
```

all the structures are jumbled up - need to align them

```r
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
```
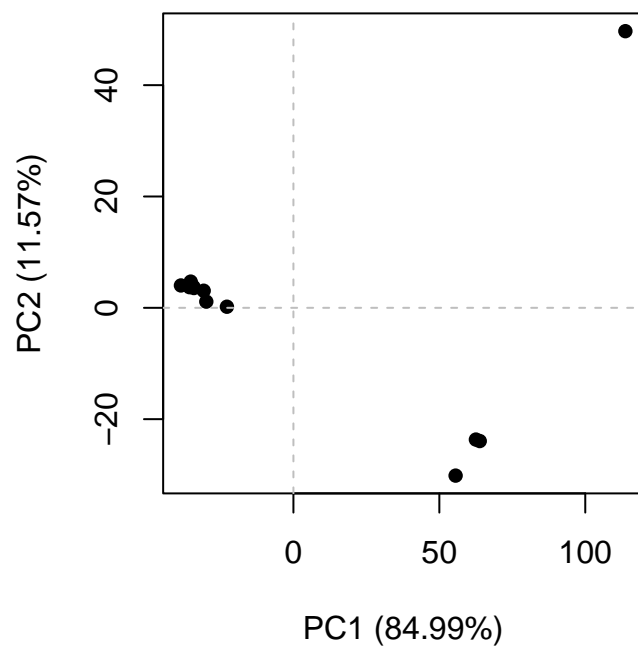
output of this is quite big, lets use PCA to visualise it

```
pc.xray<-pca(pdbs)
plot(pc.xray,1:2)
```

```
mktrj(pc.xray, file="pca_results.pdb")
```

```
#library(pak)
#view.pca(pc.xray)
```