# class09

## Mari Williams (PID: A15858833)

**PDB database**

The main repository for stuctural data for biomolecules is the Protein Data Bank (PDB) https://www.rcsb.org

```
df <- read.csv("Data Export Summary.csv")
df
```

```
          Molecular.Type   X.ray     EM    NMR Integrative Multiple.methods
1          Protein (only) 176,378 20,438 12,709         342              221
2 Protein/Oligosaccharide  10,284  3,396     34           8               11
3             Protein/NA   9,007  5,931    287          24                7
4     Nucleic acid (only)   3,077    200  1,554           2               15
5                   Other     174     13     33           3                0
6  Oligosaccharide (only)      11      0      6           0                1
  Neutron Other   Total
1      83    32 210,203
2       1     0  13,734
3       0     0  15,256
4       3     1   4,852
5       0     0     223
6       0     4      22
```

```
for (i in 2:ncol(df)) {
  df[, i] <- as.numeric(gsub(",", "", df[, i]))
}
df
```

```
          Molecular.Type  X.ray    EM   NMR Integrative Multiple.methods
1          Protein (only) 176378 20438 12709         342              221
2 Protein/Oligosaccharide  10284  3396    34           8               11
```

```
3              Protein/NA   9007  5931   287            24                   7
4      Nucleic acid (only)  3077   200  1554             2                  15
5                    Other   174    13    33             3                   0
6  Oligosaccharide (only)    11     0     6             0                   1
  Neutron Other  Total
1      83    32 210203
2       1     0  13734
3       0     0  15256
4       3     1   4852
5       0     0    223
6       0     4     22
```

This redoes the columns as numbers, but there's an easier way

```
library(readr)
```

```
Warning: package 'readr' was built under R version 4.3.3
```

```
df <- 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.
```

```
df
```

```
# A tibble: 6 x 9
  `Molecular Type`   `X-ray`    EM   NMR Integrative `Multiple methods` Neutron
  <chr>                <dbl> <dbl> <dbl>       <dbl>              <dbl>   <dbl>
1 Protein (only)      176378 20438 12709         342                221      83
2 Protein/Oligosacch~  10284  3396    34           8                 11       1
3 Protein/NA            9007  5931   287          24                  7       0
4 Nucleic acid (only)   3077   200  1554           2                 15       3
5 Other                  174    13    33           3                  0       0
6 Oligosaccharide (o~     11     0     6           0                  1       0
# i 2 more variables: Other <dbl>, Total <dbl>
```

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

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

```
[1] 81.43231
```

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

```
[1] 12.27148
```

Proportion of X-ray structures: 81.4323141 %

Proportion of EM structures: 12.2714806 %

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

```
(df[1,9] + df[2,9] + df[3,9]) / sum(df$Total) * 100
```

```
      Total
1 97.91355
```

Proportion of proteins: 97.9135453763969 %

Q3. plot something

```
library(ggplot2)
```

```
Warning: package 'ggplot2' was built under R version 4.3.3
```

```
library(dplyr)
```

```
Warning: package 'dplyr' was built under R version 4.3.3
```

```
Attaching package: 'dplyr'
```

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

    filter, lag
```
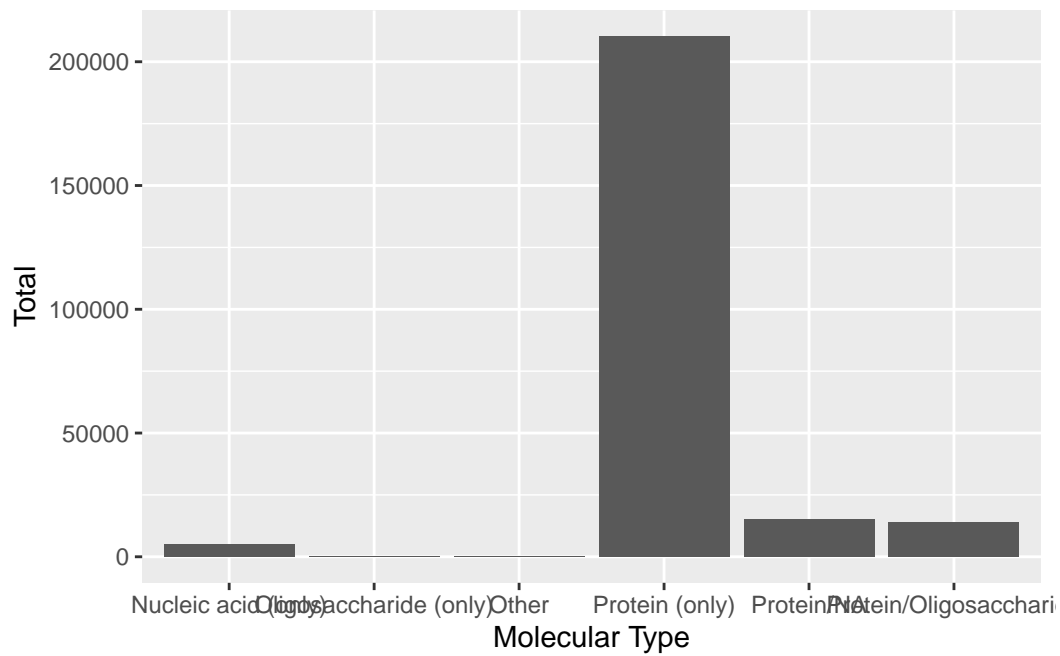
```
The following objects are masked from 'package:base':

    intersect, setdiff, setequal, union
```

```
library(tidyr)
```

```
Warning: package 'tidyr' was built under R version 4.3.3
```

```
ggplot(df) + aes(`Molecular Type`, Total) + geom_col()
```



```
df_long <- df %>%
  pivot_longer(
    cols = -`Molecular Type`,
    names_to = "Method",
    values_to = "Count"
  )%>%
  filter(Method != "Total")


df_totals <- df_long %>%
  group_by(`Molecular Type`) %>%
```
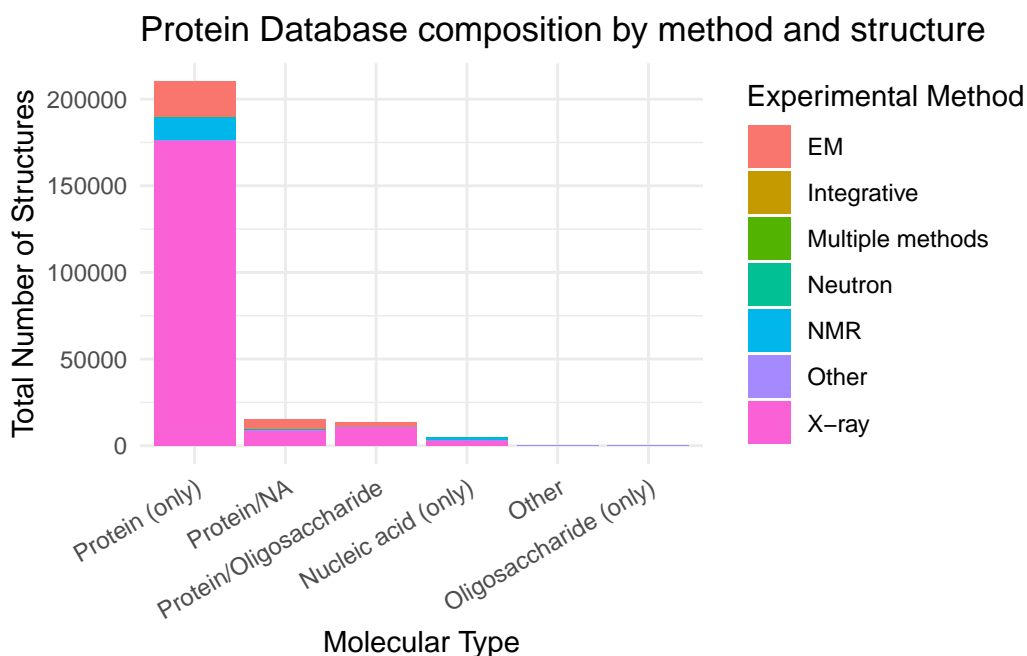
```
  summarise(total = sum(Count, na.rm = TRUE), .groups = "drop")

df_long <- df_long %>%
  left_join(df_totals, by = "Molecular Type")

ggplot(df_long, aes(x = reorder(`Molecular Type`, -total), y = Count, fill = Method)) +
  geom_bar(stat = "identity") +
  labs(
    title = "Protein Database composition by method and structure",
    x = "Molecular Type",
    y = "Total Number of Structures",
    fill = "Experimental Method"
  ) +
  theme_minimal() +
  theme(axis.text.x = element_text(angle = 30, hjust = 1))
```
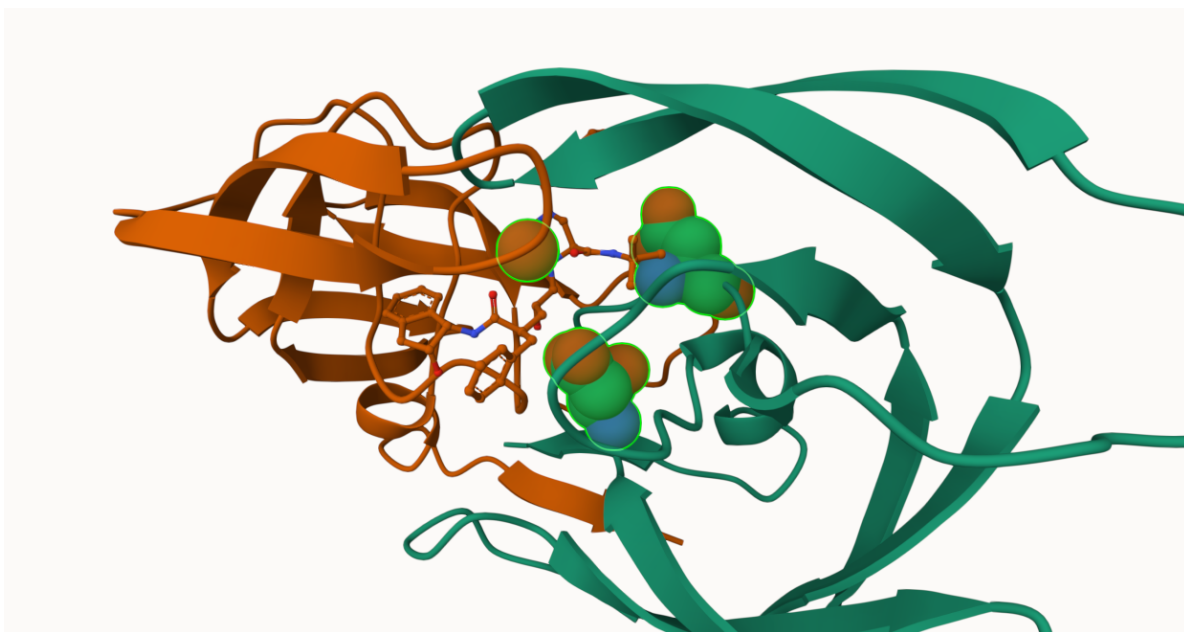


Q3: Type HIV in the PDB website search box on the home page and determine how many HIV-1 protease structures are in the current PDB?

4866

**Visualizing structure data**

The Mol* viewer is embedded into websites like Uniprot and PDB

**Bio3D package for structures**

```r
library(bio3d)
```

Warning: package 'bio3d' was built under R version 4.3.3

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

  Note: Accessing on-line PDB file

```r
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
```

```r
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
chaina <- trim.pdb(hiv, chain="A")
chainA.seq <- pdbseq(chaina)
```

```r
blast <- blast.pdb(chainA.seq)
```

```
 Searching ... please wait (updates every 5 seconds) RID = GA505F1J014
 .....
 Reporting 249 hits
```
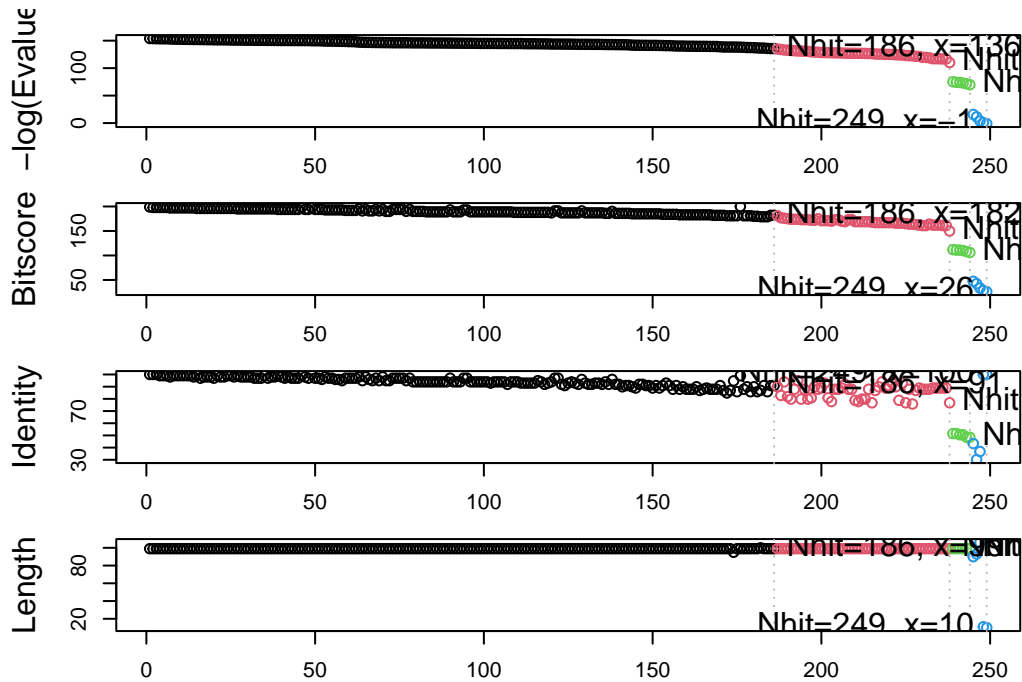
```
hits <- plot(blast)
```

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

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



## Prediction of functional motions

We can run a Normal Mode Analysis (NMA) to predict the energetic of biomolecules

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

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

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

```
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
```
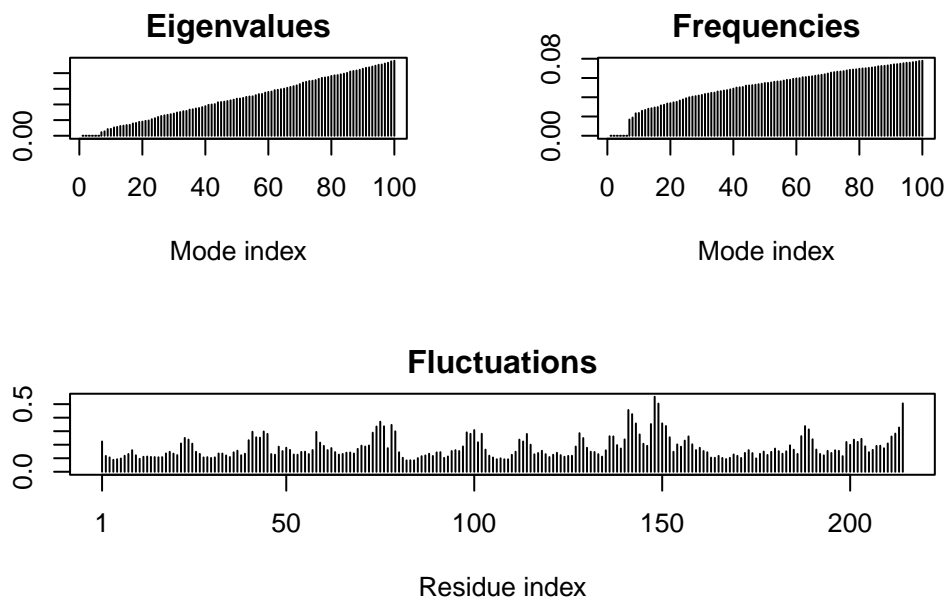
```r
m <- nma(adka)
```

```
 Building Hessian...       Done in 0.02 seconds.
 Diagonalizing Hessian...  Done in 0.18 seconds.
```

```r
plot(m)
```

**Eigenvalues**

**Frequencies**

**Fluctuations**

Residue index

Trajectory analysis

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

Install bio3dview from GitHub: install pak (install.packages('pak')) pak::pak("bioboot/bio3dview")

**Comparative analysis**

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

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

Fetching... Please wait. Done.
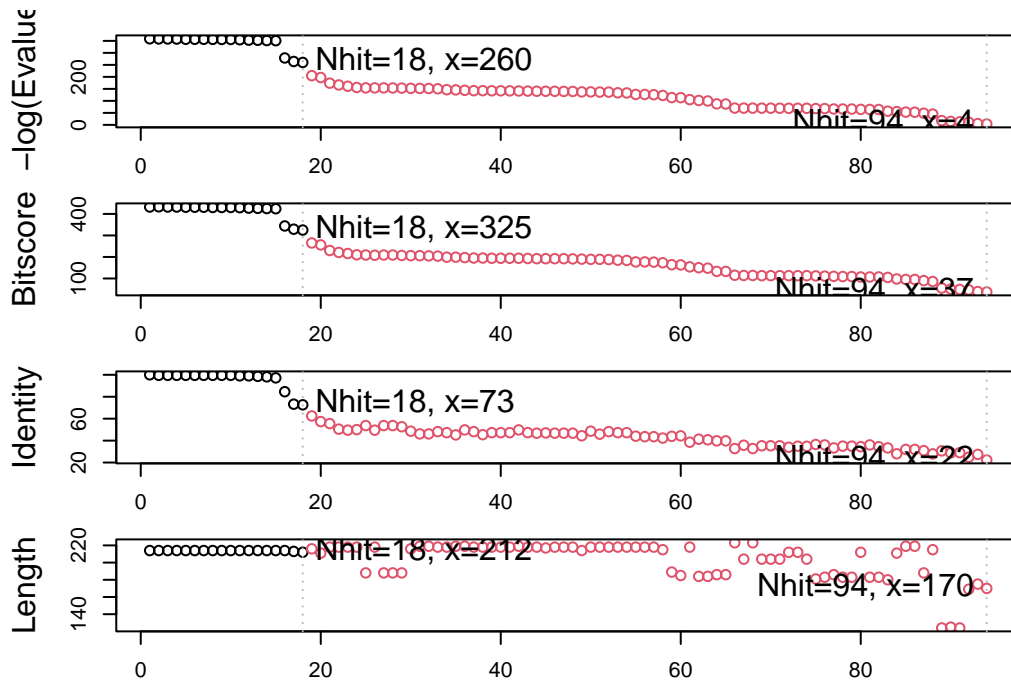
```
blast <- blast.pdb(aa)
```

 Searching ... please wait (updates every 5 seconds) RID = GX1DB088014

 Reporting 94 hits

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

```
  * Possible cutoff values:     260 3
              Yielding Nhits:    18 94

  * Chosen cutoff value of:      260
              Yielding Nhits:    18
```



```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 exists. Skipping download


Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/8BQF.pdb exists. Skipping download


Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/4X8M.pdb exists. Skipping download


Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/6S36.pdb exists. Skipping download
```

```
Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/8Q2B.pdb exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/8RJ9.pdb exists. Skipping download

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

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/4X8H.pdb exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/3HPR.pdb exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/1E4V.pdb exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/5EJE.pdb exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/1E4Y.pdb exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/3X2S.pdb exists. Skipping download

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

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

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/8PVW.pdb exists. Skipping download

Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/4K46.pdb exists. Skipping download
```

```
Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):
pdbs/4NP6.pdb exists. Skipping download
```

```
  |
  |                                                                        |   0%
  |
  |====                                                                    |   6%
  |
  |=======                                                                 |  11%
  |
  |===========                                                             |  17%
  |
  |===============                                                         |  22%
  |
  |==================                                                      |  28%
  |
  |=====================                                                   |  33%
  |
  |========================                                                |  39%
  |
  |============================                                            |  44%
  |
  |================================                                        |  50%
  |
  |====================================                                    |  56%
  |
  |========================================                                |  61%
  |
  |============================================                            |  67%
  |
  |================================================                        |  72%
  |
  |====================================================                    |  78%
  |
  |========================================================                |  83%
  |
  |============================================================            |  89%
  |
  |================================================================        |  94%
  |
  |====================================================================|  100%
```

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

Reading PDB files:
pdbs/split_chain/1AKE_A.pdb
pdbs/split_chain/8BQF_A.pdb
pdbs/split_chain/4X8M_A.pdb
pdbs/split_chain/6S36_A.pdb
pdbs/split_chain/8Q2B_A.pdb
pdbs/split_chain/8RJ9_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/8PVW_A.pdb
pdbs/split_chain/4K46_A.pdb
pdbs/split_chain/4NP6_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
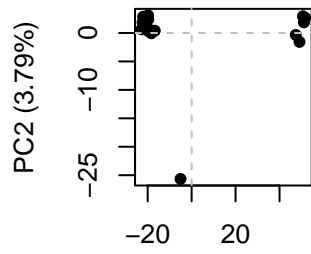..   PDB has ALT records, taking A only, rm.alt=TRUE
..   PDB has ALT records, taking A only, rm.alt=TRUE
....   PDB has ALT records, taking A only, rm.alt=TRUE
.   PDB has ALT records, taking A only, rm.alt=TRUE
.   PDB has ALT records, taking A only, rm.alt=TRUE
..


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/8BQF_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 3   name: pdbs/split_chain/4X8M_A.pdb
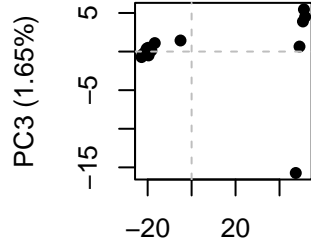pdb/seq: 4   name: pdbs/split_chain/6S36_A.pdb
```

```
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 5   name: pdbs/split_chain/8Q2B_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 6   name: pdbs/split_chain/8RJ9_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 7   name: pdbs/split_chain/6RZE_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 8   name: pdbs/split_chain/4X8H_A.pdb
pdb/seq: 9   name: pdbs/split_chain/3HPR_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 10   name: pdbs/split_chain/1E4V_A.pdb
pdb/seq: 11   name: pdbs/split_chain/5EJE_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 12   name: pdbs/split_chain/1E4Y_A.pdb
pdb/seq: 13   name: pdbs/split_chain/3X2S_A.pdb
pdb/seq: 14   name: pdbs/split_chain/6HAP_A.pdb
pdb/seq: 15   name: pdbs/split_chain/6HAM_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 16   name: pdbs/split_chain/8PVW_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 17   name: pdbs/split_chain/4K46_A.pdb
   PDB has ALT records, taking A only, rm.alt=TRUE
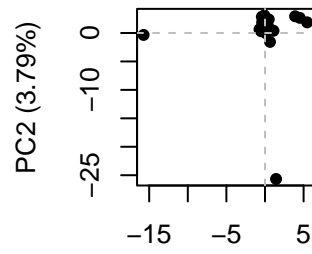pdb/seq: 18   name: pdbs/split_chain/4NP6_A.pdb
```

```r
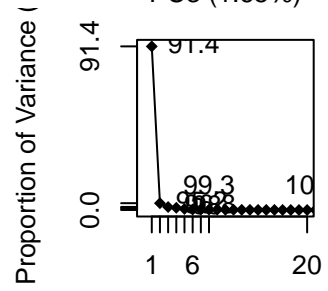pc.xray <- pca(pdbs)
plot(pc.xray)
```

By bio3dview install isnt working, sorry