

Class09

Ethan Lai

2/16/2022

First, let's look at the PDB statistics:

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

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

```
n.type <- colSums(tbl)
percentages<- (n.type/n.type["Total"] * 100)
ans<- signif(percentages, 3)
ans
```

##	X.ray	NMR	EM Multiple.methods
##	87.2000	7.2800	5.3900 0.1060
##	Neutron	Other	Total
##	0.0385	0.0198	100.0000

The proportion of Xray sructures is 87.2 % of the total The proportion of NMR sructures is 7.28 % of the total

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

```
ans2<- signif(tbl$Total[1]/ sum (tbl$Total),3) * 100
ans2
```

```
## [1] 87.3
```

87.3 % of the structures 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?

It varies quite a lot depending on our search methodology, but generally several hundred structures

Inserting a image files

```
library(bio3d)

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

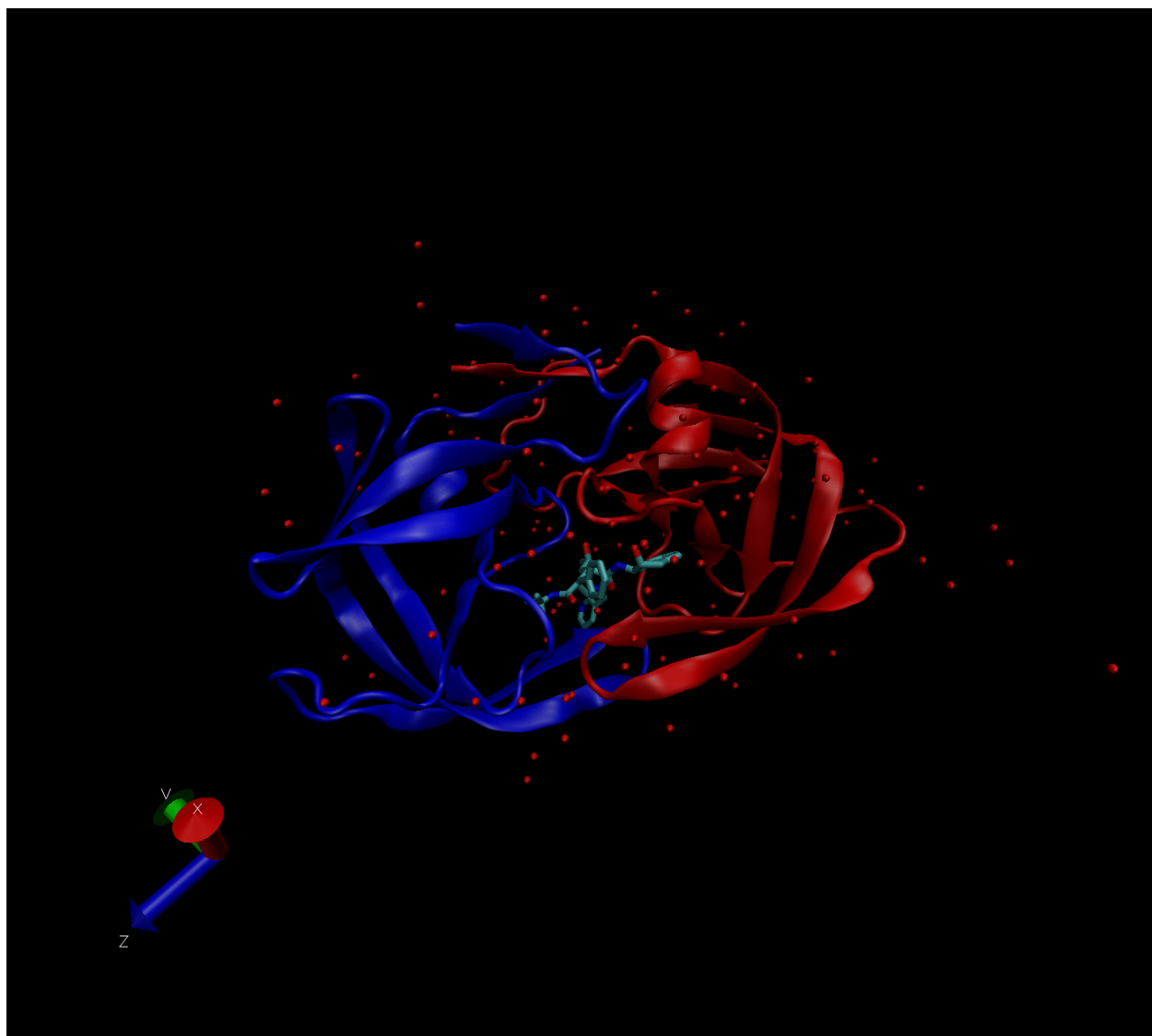


Figure 1: NS

```
## 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
```

```
aa321(c("PRO", "GLN"))
```

```
## [1] "P" "Q"
```

```
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>
```

Lets read a different single adk structure from the database now:

```
aa <- get.seq("lake_A")
```

```
## Warning in get.seq("lake_A"): Removing existing file: seqs.fasta
```

```
## Fetching... Please wait. Done.
```

```
aa
```

```
##          1      .      .      .      .      .      .      60
## pdb|1AKE|A  MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMLRAAVKSGSELGKQAKDIMDAGKLV
##          1      .      .      .      .      .      .      60
##
##          61      .      .      .      .      .      .      120
## pdb|1AKE|A  DELVIALVKERIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFDVPDELIVDRI
##          61      .      .      .      .      .      .      120
##
##          121     .      .      .      .      .      .      180
## pdb|1AKE|A  VGRRVHAPSGRVYHVKFNPPKVEGKDDVTGEELTRKDDQEETVRKRLVEYHQMTAPLIG
##          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
```

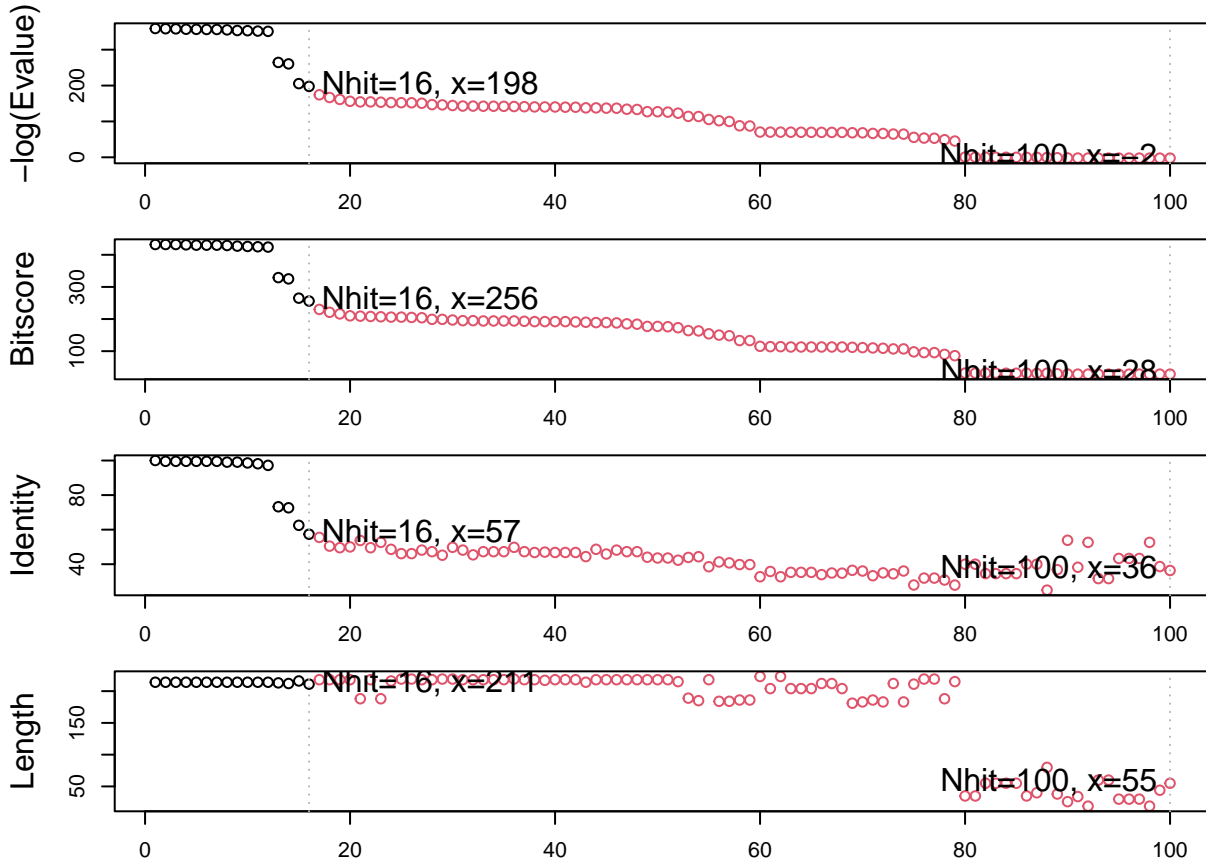
Let's find related sequences with BLAST:

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

```
## Searching ... please wait (updates every 5 seconds) RID = OV4419FG013
## .....
## Reporting 100 hits
```

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

```
## * Possible cutoff values: 197 -3
##           Yielding Nhits: 16 100
##
## * 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"
```

Now let's find an alpha fold prediction for a protein homologous to our unknown gene:

