

# JCW\_110624\_Class11-MachineLearning2

Janie Chang-Weinberg (A69037446)

Today, before delving into structure prediction with AlphaFold, we will finish off our previous lab 10 “comparative structure analysis” section.

```
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  MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMRLAAVKSGSELGKQAKDIMDAGKLV
      1      .      .      .      .      .      .      60

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

      121      .      .      .      .      .      .      180
pdb|1AKE|A  VGRRVHAPSGRVYHVKNPPKVEGKDDVTGEELTRKDDQEETVRKRLVEYHQMTAPLIG
      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
```

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

```
Searching ... please wait (updates every 5 seconds) RID = JS1PDMMP013
```

```
.
```

```
Reporting 85 hits
```

```
attributes(b)
```

```
$names
[1] "hit.tbl" "raw"      "url"
```

```
$class
[1] "blast"
```

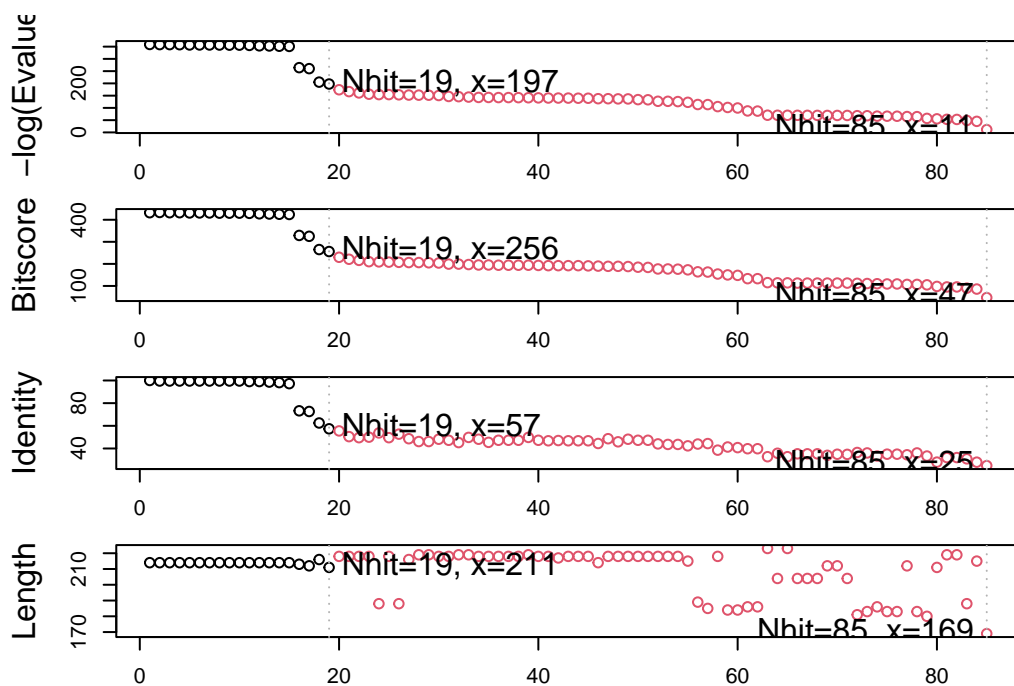
```
head(b$hit.tbl)
```

	queryid	subjectids	identity	alignmentlength	mismatches	gapopens	q.start		
1	Query_2200253	1AKE_A	100.000	214	0	0	1		
2	Query_2200253	8BQF_A	99.533	214	1	0	1		
3	Query_2200253	4X8M_A	99.533	214	1	0	1		
4	Query_2200253	6S36_A	99.533	214	1	0	1		
5	Query_2200253	8Q2B_A	99.533	214	1	0	1		
6	Query_2200253	8RJ9_A	99.533	214	1	0	1		
	q.end	s.start	s.end	evaluate	bitscore	positives	mlog.evaluate	pdb.id	acc
1	214	1	214	1.58e-156	432	100.00	358.7458	1AKE_A	1AKE_A
2	214	21	234	2.58e-156	433	100.00	358.2555	8BQF_A	8BQF_A
3	214	1	214	2.82e-156	432	100.00	358.1665	4X8M_A	4X8M_A
4	214	1	214	4.14e-156	432	100.00	357.7826	6S36_A	6S36_A
5	214	1	214	1.10e-155	431	99.53	356.8054	8Q2B_A	8Q2B_A
6	214	1	214	1.10e-155	431	99.53	356.8054	8RJ9_A	8RJ9_A

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

```
* Possible cutoff values: 197 11
    Yielding Nhits:      19 85
```

```
* Chosen cutoff value of: 197
    Yielding Nhits:      19
```



```
attributes(hits)
```

```
$names
[1] "hits"    "pdb.id"  "acc"     "inds"
```

```
$class
[1] "blast"
```

Top hits that we like from our blast results:

```
hits$pdb.id
```

```
[1] "1AKE_A" "8BQF_A" "4X8M_A" "6S36_A" "8Q2B_A" "8RJ9_A" "6RZE_A" "4X8H_A"  
[9] "3HPR_A" "1E4V_A" "5EJE_A" "1E4Y_A" "3X2S_A" "6HAP_A" "6HAM_A" "4K46_A"  
[17] "4NP6_A" "3GMT_A" "4PZL_A"
```

```
#Download related PDB files  
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/8BQF.pdb.gz exists. Skipping download
```

```
Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):  
pdbs/4X8M.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/8Q2B.pdb.gz exists. Skipping download
```

```
Warning in get.pdb(hits$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):  
pdbs/8RJ9.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/4X8H.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

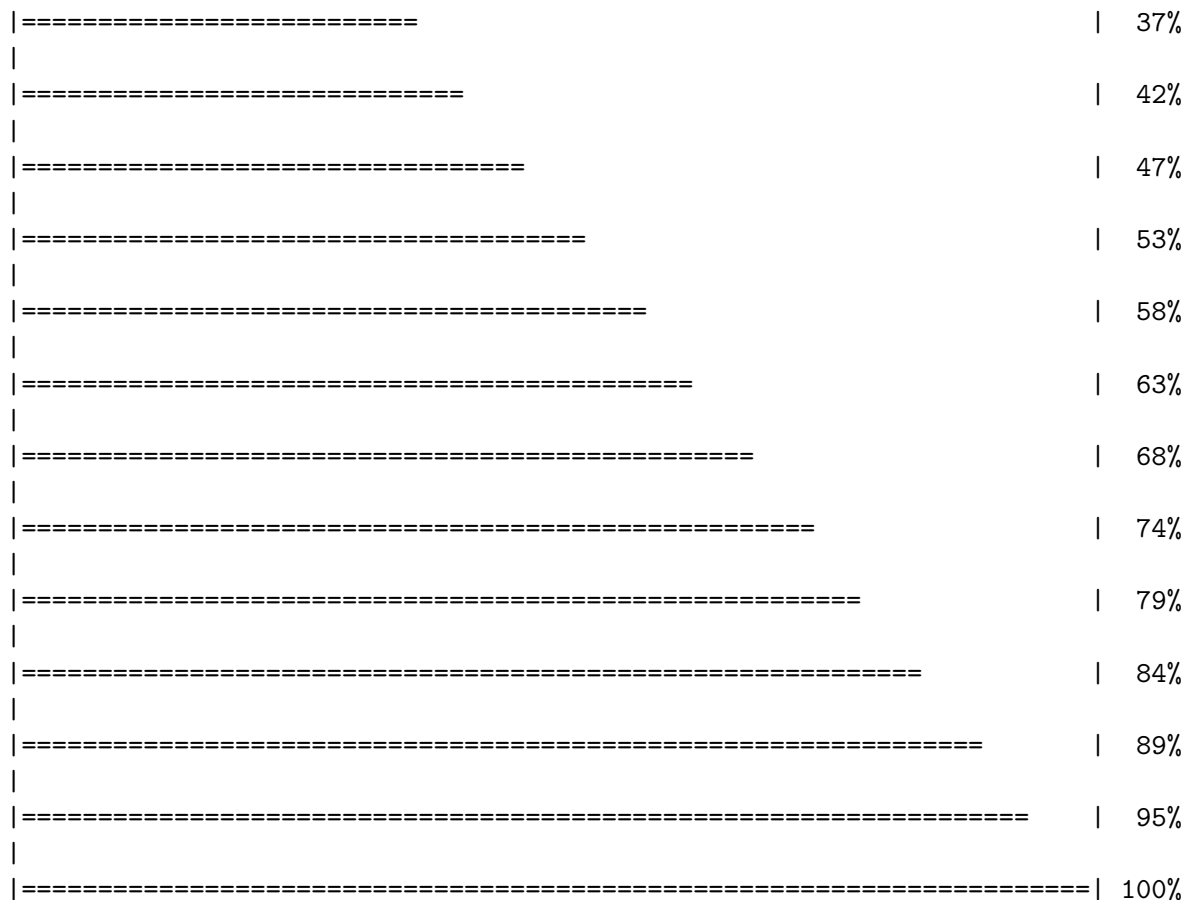
Warning in get.pdb(hits\$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):  
pdbs/4K46.pdb.gz exists. Skipping download

Warning in get.pdb(hits\$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):  
pdbs/4NP6.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%
====	5%
=====	11%
=====	16%
=====	21%
=====	26%
=====	32%



I have now found and downloaded all ADK structures in the PDB database but viewing them is difficult as they need to be aligned and superposed.

I am going to install the BiocManager from CRAN, then I can use `BiocManager::install()` to install any bioconductor package.

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

Reading PDB files:

```
pdbbs/split_chain/1AKE_A.pdb
pdbbs/split_chain/8BQF_A.pdb
pdbbs/split_chain/4X8M_A.pdb
pdbbs/split_chain/6S36_A.pdb
pdbbs/split_chain/8Q2B_A.pdb
pdbbs/split_chain/8RJ9_A.pdb
pdbbs/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/4K46_A.pdb
pdbs/split_chain/4NP6_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
..  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/4K46_A.pdb
             PDB has ALT records, taking A only, rm.alt=TRUE
pdb/seq: 17  name: pdbs/split_chain/4NP6_A.pdb
pdb/seq: 18  name: pdbs/split_chain/3GMT_A.pdb
pdb/seq: 19  name: pdbs/split_chain/4PZL_A.pdb

```

## pdbs

```

[Truncated_Name:1] 1AKE_A.pdb      1      .      .      .      40
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:2] 8BQF_A.pdb      1      .      .      .      40
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:3] 4X8M_A.pdb      1      .      .      .      40
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:4] 6S36_A.pdb      1      .      .      .      40
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:5] 8Q2B_A.pdb      1      .      .      .      40
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:6] 8RJ9_A.pdb      1      .      .      .      40
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:7] 6RZE_A.pdb      1      .      .      .      40
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:8] 4X8H_A.pdb      1      .      .      .      40
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:9] 3HPR_A.pdb      1      .      .      .      40
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:10] 1E4V_A.pdb     1      .      .      .      40
-----MRIILLGAPVAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:11] 5EJE_A.pdb     1      .      .      .      40
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:12] 1E4Y_A.pdb     1      .      .      .      40
-----MRIILLGALVAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:13] 3X2S_A.pdb     1      .      .      .      40
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:14] 6HAP_A.pdb     1      .      .      .      40
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:15] 6HAM_A.pdb     1      .      .      .      40
-----MRIILLGAPGAGKGTQAQFIMEKYGIPQIS
[Truncated_Name:16] 4K46_A.pdb     1      .      .      .      40
-----MRIILLGAPGAGKGTQAQFIMAKFGIPQIS
[Truncated_Name:17] 4NP6_A.pdb     1      .      .      .      40
-----NAMRIILLGAPGAGKGTQAQFIMEKFGIPQIS
[Truncated_Name:18] 3GMT_A.pdb     1      .      .      .      40
-----MRLILLGAPGAGKGTQANFIKEKFGIPQIS
[Truncated_Name:19] 4PZL_A.pdb     1      .      .      .      40
TENLYFQSNAMRIILLGAPGAGKGTQAKIIEQKYNIAHIS
             **^*****  *****  *  *^ *  **
1      .      .      .      40

41      .      .      .      80
TGDMLRAAVKSGSELGKQAKDIMDAGKLVTDLVIALVKE
[Truncated_Name:1] 1AKE_A.pdb      TGDMLRAAVKSGSELGKQAKDIMDAGKLVTDLVIALVKE
[Truncated_Name:2] 8BQF_A.pdb      TGDMLRAAVKSGSELGKQAKDIMDAGKLVTDLVIALVKE
[Truncated_Name:3] 4X8M_A.pdb      TGDMLRAAVKSGSELGKQAKDIMDAGKLVTDLVIALVKE

```



[Truncated_Name:4] 6S36_A.pdb	TGDMRLRAAVKSGSELGKQAKDIMDAGKLVTDDELVIALVKE
[Truncated_Name:5] 8Q2B_A.pdb	TGDMRLRAAVKSGSELGKQAKDIMDAGKLVTDDELVIALVKE
[Truncated_Name:6] 8RJ9_A.pdb	TGDMRLRAAVKSGSELGKQAKDIMDAGKLVTDDELVIALVKE
[Truncated_Name:7] 6RZE_A.pdb	TGDMRLRAAVKSGSELGKQAKDIMDAGKLVTDDELVIALVKE
[Truncated_Name:8] 4X8H_A.pdb	TGDMRLRAAVKSGSELGKQAKDIMDAGKLVTDDELVIALVKE
[Truncated_Name:9] 3HPR_A.pdb	TGDMRLRAAVKSGSELGKQAKDIMDAGKLVTDDELVIALVKE
[Truncated_Name:10] 1E4V_A.pdb	TGDMRLRAAVKSGSELGKQAKDIMDAGKLVTDDELVIALVKE
[Truncated_Name:11] 5EJE_A.pdb	TGDMRLRAAVKSGSELGKQAKDIMDACKLVTDDELVIALVKE
[Truncated_Name:12] 1E4Y_A.pdb	TGDMRLRAAVKSGSELGKQAKDIMDAGKLVTDDELVIALVKE
[Truncated_Name:13] 3X2S_A.pdb	TGDMRLRAAVKSGSELGKQAKDIMDCGKLVTDDELVIALVKE
[Truncated_Name:14] 6HAP_A.pdb	TGDMRLRAAVKSGSELGKQAKDIMDAGKLVTDDELVIALVRE
[Truncated_Name:15] 6HAM_A.pdb	TGDMRLRAAIKSGSELGKQAKDIMDAGKLVTDDEIIIALVKE
[Truncated_Name:16] 4K46_A.pdb	TGDMRLRAAIKAGTELGKQAKSVIDAGQLVSDDIILGLVKE
[Truncated_Name:17] 4NP6_A.pdb	TGDMRLRAAIKAGTELGKQAKAVIDAGQLVSDDIILGLIKE
[Truncated_Name:18] 3GMT_A.pdb	TGDMRLRAAVKAGTPLGVEAKTYMDEGKLVPSLIIGLVKE
[Truncated_Name:19] 4PZL_A.pdb	TGDMIRETIKSGSALGQELKKVLDAGELVSDEFIIVKIVKD
	****~* ~* *~ ** * ~* ** * ~ ~~~~~
	41 . . . 80
	81 . . . 120
[Truncated_Name:1] 1AKE_A.pdb	RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:2] 8BQF_A.pdb	RIAQE----GFLLDGFPR TIPQADAMKEAGINVDYVIEFD
[Truncated_Name:3] 4X8M_A.pdb	RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:4] 6S36_A.pdb	RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:5] 8Q2B_A.pdb	RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:6] 8RJ9_A.pdb	RIAQEDCRNGFLLAGFPRTIPQADAMKEAGINVDYVLEFD
[Truncated_Name:7] 6RZE_A.pdb	RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:8] 4X8H_A.pdb	RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:9] 3HPR_A.pdb	RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:10] 1E4V_A.pdb	RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:11] 5EJE_A.pdb	RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:12] 1E4Y_A.pdb	RIAQEDCRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:13] 3X2S_A.pdb	RIAQEDSRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:14] 6HAP_A.pdb	RICQEDSRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:15] 6HAM_A.pdb	RICQEDSRNGFLLDGFPR TIPQADAMKEAGINVDYVLEFD
[Truncated_Name:16] 4K46_A.pdb	RIAQDDCAKGFLLDGFPR TIPQADGLKEVGVVVDYVIEFD
[Truncated_Name:17] 4NP6_A.pdb	RIAQADCEKGFLLDGFPR TIPQADGLKEMGINVDYVIEFD
[Truncated_Name:18] 3GMT_A.pdb	RLKEADCANGYLFDFPRTIAQADAMKEAGVAIDYVLEID
[Truncated_Name:19] 4PZL_A.pdb	RISKNDCCNNGFLLDGVPR TIPQAQELDKLVNIDYIVEVD
	*~ *~* * **** ** ^ *~ ~**~* *
	81 . . . 120
	121 . . . 160



161

200

201

227

```

[Truncated_Name:1] 1AKE_A.pdb      T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:2] 8BQF_A.pdb      T--KYAKVDGTKPVAEVRADLEKIL--
[Truncated_Name:3] 4X8M_A.pdb      T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:4] 6S36_A.pdb      T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:5] 8Q2B_A.pdb      T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:6] 8RJ9_A.pdb      T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:7] 6RZE_A.pdb      T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:8] 4X8H_A.pdb      T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:9] 3HPR_A.pdb      T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:10] 1E4V_A.pdb     T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:11] 5EJE_A.pdb     T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:12] 1E4Y_A.pdb     T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:13] 3X2S_A.pdb     T--KYAKVDGTKPVAEVRADLEKILG-
[Truncated_Name:14] 6HAP_A.pdb     T--KYAKVDGTKPVCEVRADLEKILG-
[Truncated_Name:15] 6HAM_A.pdb     T--KYAKVDGTKPVCEVRADLEKILG-
[Truncated_Name:16] 4K46_A.pdb     T--QYLKFDGTKAVAEVSAELEKALA-
[Truncated_Name:17] 4NP6_A.pdb     T--QYLKFDGTKQVSEVSADIAKALA-
[Truncated_Name:18] 3GMT_A.pdb     E-----NGLKAPA-----YRKISG-
[Truncated_Name:19] 4PZL_A.pdb     KIPKYIKINGDQAVEKVSQDIFDQLNK

```

\*

201

227

Call:

```
pdbaln(files = files, fit = TRUE, exeFile = "msa")
```

Class:

```
pdbs, fasta
```

Alignment dimensions:

```
19 sequence rows; 227 position columns (199 non-gap, 28 gap)
```

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

```
##Principal Component Analysis
```

```
pc <- pca(pdb)
plot(pc)
```

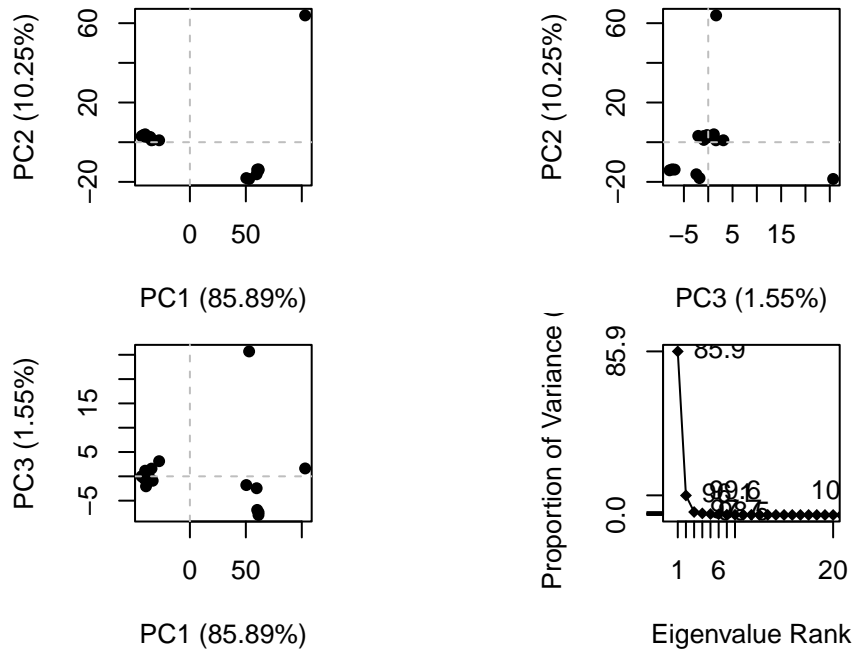
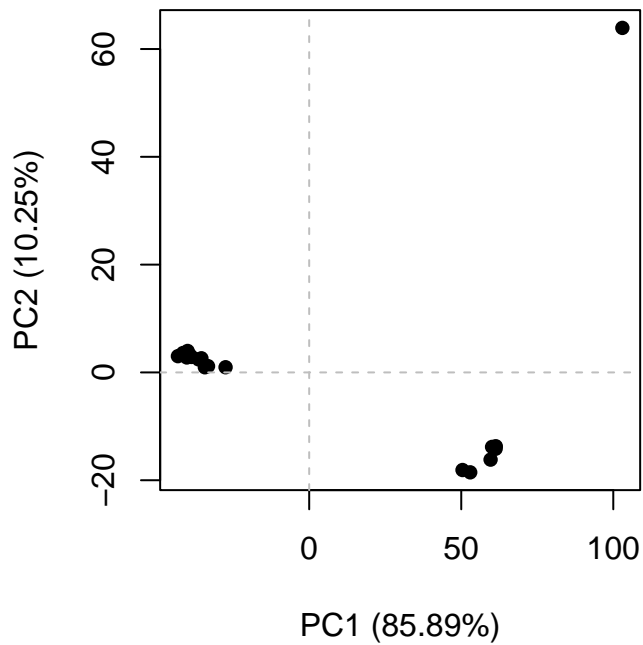


fig 4 here is the scree plot! “Eigenvalue” are how you order the PCs. Eigenvalue 1 is PC1, Eigenvalue 2 is PC2, etc.

```
plot(pc, pc.axes=c(1:2))
```



To examine in more detail what PC1 (or any PC) is capturing here, we can plot the loadings or make an animated file of moving along PC1

```
mktrj(pc, pc=1, file="pc1.pdb")
```