```
Class 9 Lab
AUTHOR
Catherine Diep
The RCSB Protein Data Bank (PDB)
Protein structures by X-ray crystalgraphy dominate this database. We are skipping Q1-3 as the
website was too slow for us.
2. Visualizing the HIV-1 protease structure
HIV-Pr structure from 1hsg
  Q4: Water molecules normally have 3 atoms. Why do we see just one atom per water
  molecule in this structure?
We only see one atom (the oxygen atom) per molecule because the hydrogen atoms are too small
to be displayed by PDB.
  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?
This water is labeled HOH 308.
  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 (we recommend "Ball & Stick" for these side-chains). Add this figure to your Quarto
  document.
1hsg with marked Asp
  Discussion Topic: Can you think of a way in which indinavir, or even larger ligands and
  substrates, could enter the binding site? Larger ligands could enter the binding site when
  the HIV-protease's two chains move apart to expose the site.
3. Introduction to Bio3D in R.
Bio3D is an R package for structural bioinformatics. To use it we need to call it up iwth the 'libary()'
function.
 library(bio3d)
 pdb <- read.pdb("1hsg")</pre>
  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
  Q7: How many amino acid residues are there in this pdb object? 198
  Q8: Name one of the two non-protein residues? HOH
  Q9: How many protein chains are in this structure? 2
The ATAM records of a PDB file are stored in 'pdb$atom'.
 head(pdb$atom)
  type eleno elety alt resid chain resno insert
                                                                      Z 0
                                              <NA> 29.361 39.686 5.862 1 38.10
1 ATOM
                 N < NA >
                           PR0
                                              <NA> 30.307 38.663 5.319 1 40.62
2 ATOM
                CA <NA>
                           PR0
                                              <NA> 29.760 38.071 4.022 1 42.64
3 ATOM
                 C <NA>
                           PR0
                                              <NA> 28.600 38.302 3.676 1 43.40
4 ATOM
                  0 <NA>
                           PR0
                                              <NA> 30.508 37.541 6.342 1 37.87
5 ATOM
                CB <NA>
                           PR0
6 ATOM
                 CG <NA>
                                              <NA> 29.296 37.591 7.162 1 38.40
                           PR0
  segid elesy charge
1 <NA>
                 < NA>
   <NA>
                 < NA>
   <NA>
                 <NA>
   <NA>
                 < NA>
   <NA>
                 < NA>
   <NA>
                 <NA>
Comparative analysis of Adenylate kinase
(ADK)
  Q10. Which of the packages above is found only on BioConductor and not CRAN? MSA
  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
We will start our analysis with a single PDB id (code from the PDB database): 1AKE
First we get its primary sequence:
 aa <- get.seq("1ake_a")</pre>
Warning in get.seq("1ake_a"): Removing existing file: seqs.fasta
Fetching... Please wait. Done.
 aa
                                                                           60
pdb|1AKE|A MRIILLGAPGAGKGTQAQFIMEKYGIPQISTGDMLRAAVKSGSELGKQAKDIMDAGKLVT
                                                                           120
             61
pdb|1AKE|A DELVIALVKERIAQEDCRNGFLLDGFPRTIPQADAMKEAGINVDYVLEFDVPDELIVDRI
                                                                           120
           121
                                                                           180
pdb|1AKE|A VGRRVHAPSGRVYHVKFNPPKVEGKDDVTGEELTTRKDDQEETVRKRLVEYHQMTAPLIG
           121
                                                                           180
           181
                                                214
pdb|1AKE|A YYSKEAEAGNTKYAKVDGTKPVAEVRADLEKILG
           181
Call:
  read.fasta(file = outfile)
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? 214
 # Blast or hmmr search
 #b <- blast.pdb(aa)</pre>
 #hits <- plot(b)</pre>
 #List out some 'top hits'
 #head(hits$pdb.id)
Use these ADK structures for analysis:
 hits <- NULL
 hits$pdb.id <- c('1AKE_A','6S36_A','6RZE_A','3HPR_A','1E4V_A','5EJE_A','1E4Y_A','
Download all these PDB files from the database...
 files <- get.pdb(hits$pdb.id, path="pdbs", split=TRUE, gzip=TRUE)</pre>
Warning in get.pdb(hits$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%
```

Align releated PDBs

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

Reading PDB files:

pdbs <- pdbaln(files, fit = TRUE, exefile="msa")</pre>

23%

31%

38%

46%

54%

62%

69%

77%

85%

92%

pdbs/split_chain/3GMT_A.pdb pdbs/split_chain/4PZL_A.pdb 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 name: pdbs/split_chain/4PZL_A.pdb pdb/seq: 13 # Vector containing PDB codes for figure axis ids <- basename.pdb(pdbs\$id)</pre> # Draw schematic alignment plot(pdbs, labels=ids) Sequence Alignment Overview 6S36_A 1AKE_A 6RZE_A

3HPR_A

1E4Y_A

1E4V_A

5EJE_A

3X2S_A

6HAM_A

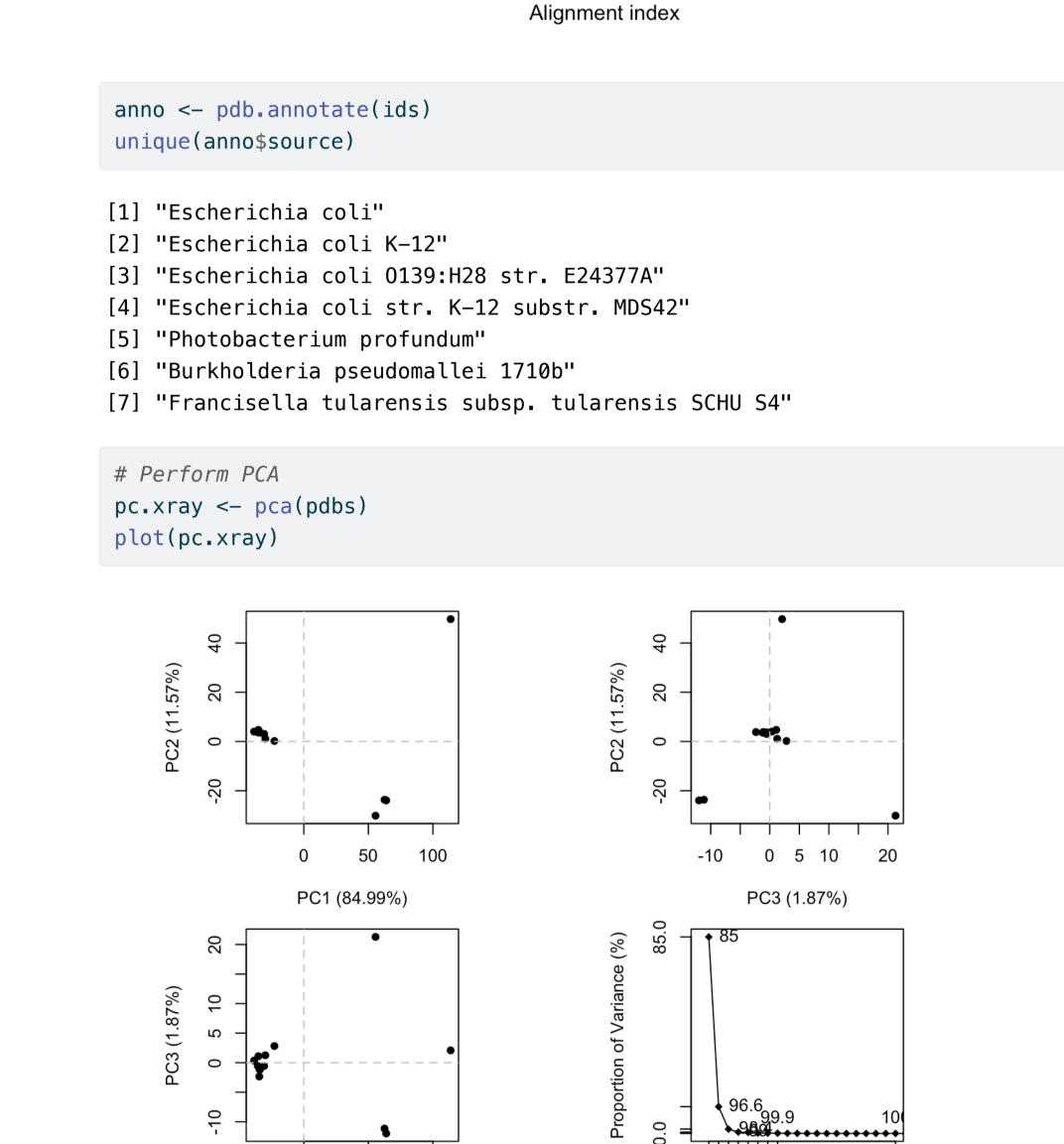
6HAP_A

4K46_A

3GMT_A

4PZL_A

200



50

PC1 (84.99%)

Calculate RMSD

rd <- rmsd(pdbs)</pre>

Structure-based clustering

grps.rd <- cutree(hc.rd, k=3)</pre>

hc.rd <- hclust(dist(rd))</pre>

PC2

0

-20

100

100

150

50

plot(pc.xray, 1:2, col="grey50", bg=grps.rd, pch=21, cex=1) 40 (11.57%)20

50

PC1 (84.99%)

0

100

Warning in rmsd(pdbs): No indices provided, using the 204 non NA positions

1 4 7

Eigenvalue Rank

20