

# Class09

Nicholas Thiphakhinkeo A17686679

## Read File

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

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

```
pdbdb$Total
```

```
[1] "195,866" "12,328" "13,746" "4,532" "213" "22"
```

## Remove Commas Convert to Numerics

```
as.numeric( sub(",", "", pdbdb$Total) )
```

```
[1] 195866 12328 13746 4532 213 22
```

## Set as Function for Whole Table and Future Tables

```
x <- pdbdb$Total  
as.numeric( sub(",", "", x))
```

```
[1] 195866 12328 13746 4532 213 22
```

```
comma2numeric <- function(x) {  
  as.numeric( sub(",", "", x))  
}
```

## Test

```
comma2numeric(pdbdb$X.ray)
```

```
[1] 167317 9645 8735 2869 170 11
```

```
apply(pdbdb, 2, comma2numeric)
```

Warning in FUN(newX[, i], ...): NAs introduced by coercion

	Molecular.Type	X.ray	EM	NMR	Multiple.methods	Neutron	Other	Total
[1,]	NA	167317	15698	12534	208	77	32	195866
[2,]	NA	9645	2639	34	8	2	0	12328
[3,]	NA	8735	4718	286	7	0	0	13746
[4,]	NA	2869	138	1507	14	3	1	4532
[5,]	NA	170	10	33	0	0	0	213
[6,]	NA	11	0	6	1	0	4	22

## Try different Read/Import Function

```
library(readr)  
pdbdb <- read_csv("Data Export Summary.csv")
```

Rows: 6 Columns: 8

-- Column specification -----

Delimiter: ","

chr (1): Molecular Type

dbl (3): 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.

```
sum(pdbdb$Total)
```

```
[1] 226707
```

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

```
[1] 83.25592
```

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

```
[1] 10.2348
```

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

```
pdbdb$Total[1]/ sum(pdbdb$Total) * 100
```

```
[1] 86.3961
```

## 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?

Bio3D

```
library(bio3d)
pdb <- read.pdb("1hsg")
```

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?

```
sum(pdb$calpha)
```

```
[1] 198
```

```
length(pdbseq(pdb))
```

```
[1] 198
```

Q8: Name one of the two non-protein residues

HOH, MK1

Q9: How many protein chains are in this structure?

```
attributes(pdb)
```

```
$names
```

```
[1] "atom" "xyz" "seqres" "helix" "sheet" "calpha" "remark" "call"
```

```
$class
```

```
[1] "pdb" "sse"
```

```
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>

```
adk <- read.pdb("6s36")
```

Note: Accessing on-line PDB file

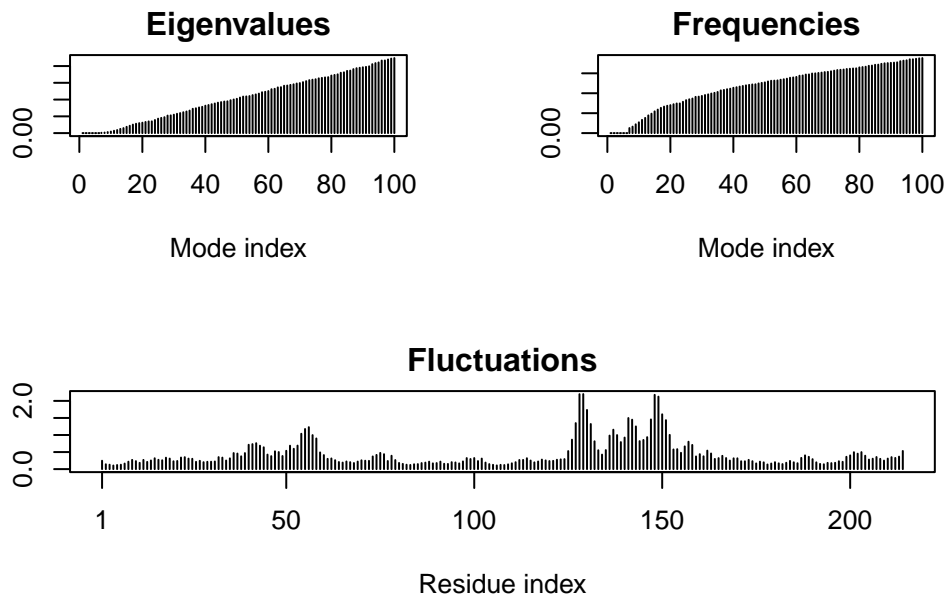
PDB has ALT records, taking A only, rm.alt=TRUE

## Flexibility Prediction

```
m <- nma(adk)
```

```
Building Hessian...      Done in 0.05 seconds.  
Diagonalizing Hessian... Done in 0.36 seconds.
```

```
plot(m)
```



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

MSA Package

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

Q12. T/F? Functions from the devtools package can be used to install packages from GitHub and BitBucket?

True

```
library(bio3d)  
aa <- get.seq("lake_A")
```

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

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

Q13. How many amino acids are in this sequence, i.e. how long is this sequence?

## Blast Search

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

## Plot summar of Search Results

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

## List out some 'top hits'

```
# head(hits$ pdb.id)
```

```
hits <- NULL
```

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

## Download releated PDB files

```
files <- get.pdb(hits$ pdb.id, path="pdb", split=TRUE, gzip=TRUE)
```

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

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

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

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

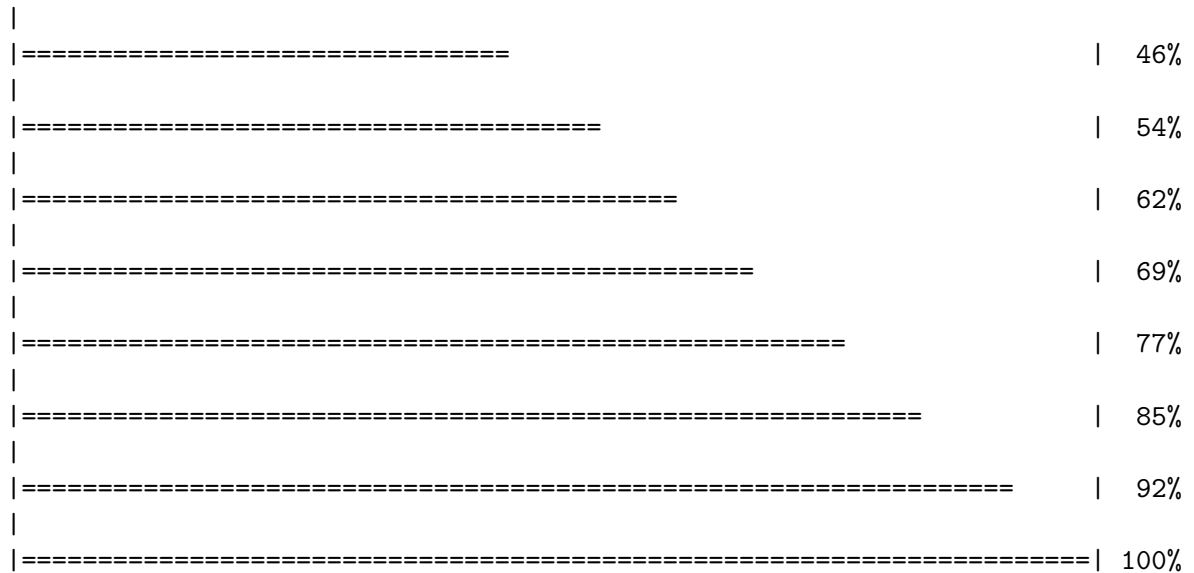
Warning in get.pdb(hits\$pdb.id, path = "pdbs", split = TRUE, gzip = TRUE):  
pdbs/4K46.pdb exists. Skipping download

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

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

	0%
=====	8%
=====	15%
=====	23%
=====	31%
=====	38%





## Align related PDBs

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

Reading PDB files:

pdbbs/split\_chain/1AKE\_A.pdb

pdbbs/split\_chain/6S36\_A.pdb

pdbbs/split\_chain/6RZE\_A.pdb

pdbbs/split\_chain/3HPR\_A.pdb

pdbbs/split\_chain/1E4V\_A.pdb

pdbbs/split\_chain/5EJE\_A.pdb

pdbbs/split\_chain/1E4Y\_A.pdb

pdbbs/split\_chain/3X2S\_A.pdb

pdbbs/split\_chain/6HAP\_A.pdb

pdbbs/split\_chain/6HAM\_A.pdb

pdbbs/split\_chain/4K46\_A.pdb

pdbbs/split\_chain/3GMT\_A.pdb

pdbbs/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
```

## Vector containing PDB codes for figure axis

```
ids <- basename.pdb(pdb$id)
```

## Draw schematic alignment

```
#plot(pdb, labels=ids)
```

## Annotating PDB Structures

```
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"
```

```
anno
```

	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 (ADK)		
3HPR_A	2.00	<NA>	Adenylate kinase, active site lid (ADK_lid)		
1E4V_A	1.85	Adenylate kinase	Adenylate kinase (ADK)		
5EJE_A	1.90	<NA>	Adenylate kinase, active site lid (ADK_lid)		
1E4Y_A	1.85	Adenylate kinase	Adenylate kinase, active site lid (ADK_lid)		
3X2S_A	2.80	<NA>	Adenylate kinase (ADK)		
6HAP_A	2.70	<NA>	Adenylate kinase (ADK)		
6HAM_A	2.55	<NA>	Adenylate kinase, active site lid (ADK_lid)		

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 (ADK)

	ligandId
1AKE_A	AP5
6S36_A	CL (3),NA,MG (2)
6RZE_A	CL (2),NA (3)
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	CHLORIDE ION (2),SODIUM ION (3)
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 0139:H28 str. E24377A
1E4Y_A	Escherichia coli
3X2S_A	Escherichia coli str. K-12 substr. MDS42
6HAP_A	Escherichia coli 0139: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  
 1E4Y\_A  
 3X2S\_A  
 6HAP\_A  
 6HAM\_A  
 4K46\_A  
 3GMT\_A  
 4PZL\_A

Cryst

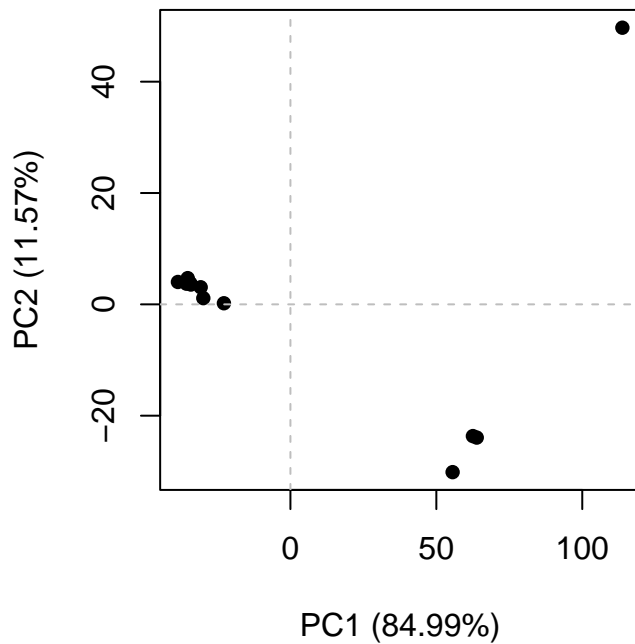
The crys

		citation	rObserved	rFree
1AKE_A		Muller, C.W., et al. J Mol Biol (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

## Perform PCA

```
pc.xray <- pca(pdbbs)
plot(pc.xray, pc.axes = c(1,2))
```



## Calculate RMSD

```
rd <- rmsd(pdbbs)
```

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

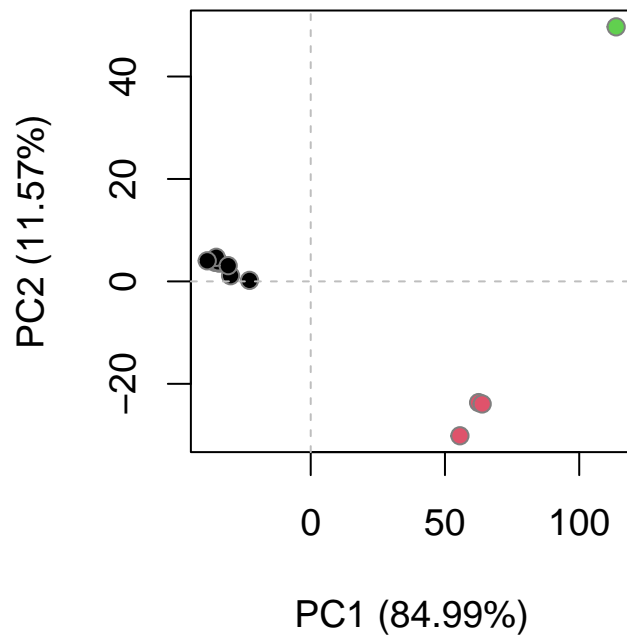
## Structure-based Clustering

```

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)

```



## Visualize first principal component

```
pc1 <- mktrj(pc.xray, pc=1, file="pc_1.pdb")
```

#Plotting results with ggplot2

```

library(ggplot2)
library(ggrepel)

df <- data.frame(PC1=pc.xray$z[,1],
                  PC2=pc.xray$z[,2],
                  col=as.factor(grps.rd),
                  ids=ids)

```

```
p <- ggplot(df) +
  aes(PC1, PC2, col=col, label=ids) +
  geom_point(size=2) +
  geom_text_repel(max.overlaps = 20) +
  theme(legend.position = "none")
p
```

