

# Data analysis challenge

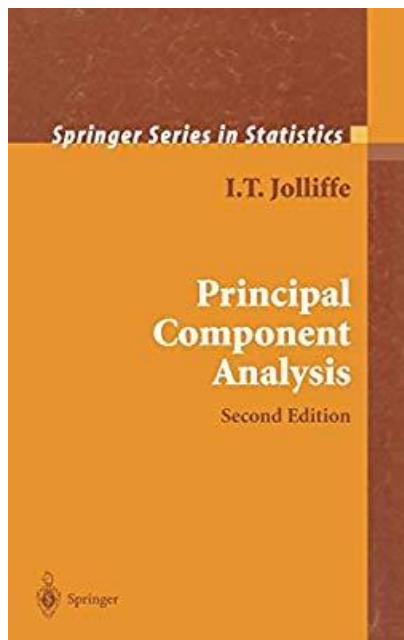
- Aim: to understand what is «wrong» with the mutated protein
- Why? We need a way of comparing movements among the different molecular simulations.
- How?
  - Because this analysis requires a lot of processing, you will have access to Power9 at BSC.
  - 3 main steps: generate simplified trajectories, perform analysis, and visualize the data.
  - Example scripts are provided for 2 analyses:
    - Principal Component Analysis (PCA) to see the most dominant motions (with most variance)
    - Autocorrelation Map to see correlated or anti-correlated motions
- Questions:
  - Is it possible to parallelize the calculations?
  - How can PCA be used to:
    - compare the movement of different chlorine and protonation states («healthy» or mutant protein)?
    - compare the movements of the «healthy» with the mutated protein?
  - Can you find a graphical way to easily compare 2 autocorrelation maps? e.g.
    - Comparison of the same chlorine state but a different protonation state for the «healthy» protein
    - Comparison of the same chlorine state but a different protonation state for the mutant protein
    - Comparison of the same chlorine and protonation state of «healthy» protein vs. mutant

## Could PCA be improved in this analysis of simulations?

The loop which has larger contribution in the PCA is also located farthest from the center of coordinates (in the center of protein) and PCA analysis doesn't show other regions significantly contributing.

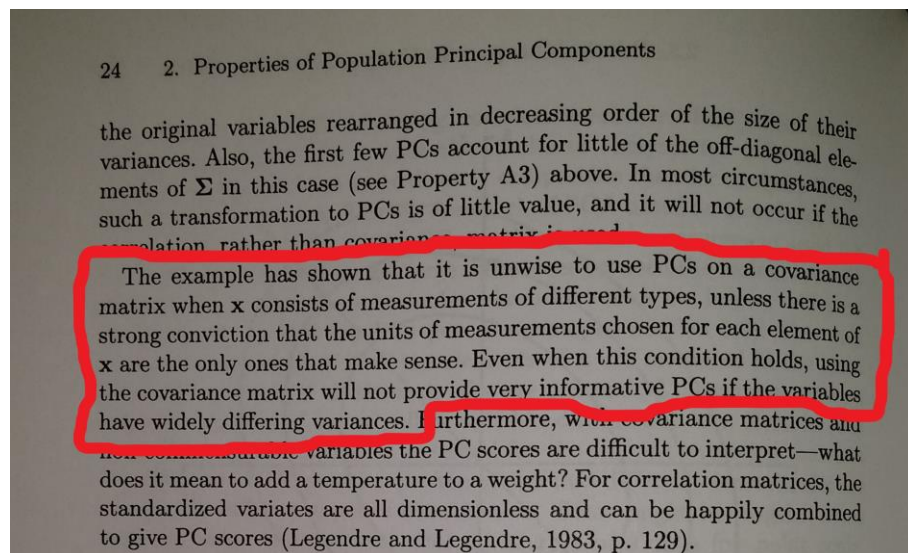
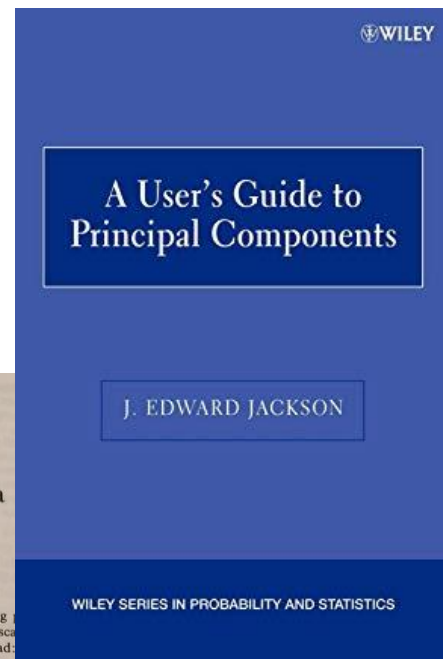


# PCA textbooks recommend scaling & centering when variances are different between variables



For instance:

$$x'_i = \frac{(x_i - \bar{x})}{std. dev_i}$$



## When they are not centered and scaled, fewer PC dominate and mask the effect of other PCs

Dataset of cars (included by default in R package)

```
> summary(mtcars)
```

mpg	cyl
Min. :10.40	Min. :4.000
1st Qu.:15.43	1st Qu.:4.000
Median :19.20	Median :6.000
Mean :20.09	Mean :6.188
3rd Qu.:22.80	3rd Qu.:8.000
Max. :33.90	Max. :8.000

drat	wt
Min. :2.760	Min. :1.513
1st Qu.:3.080	1st Qu.:2.581
Median :3.695	Median :3.325
Mean :3.597	Mean :3.217
3rd Qu.:3.920	3rd Qu.:3.610
Max. :4.930	Max. :5.424

am
Min. :0.0000
1st Qu.:0.0000
Median :0.0000
Mean :0.4062
3rd Qu.:1.0000
Max. :1.0000

gear
Min. :3.000
1st Qu.:3.000
Median :4.000
Mean :3.688
3rd Qu.:4.000
Max. :5.000

disp	hp
Min. : 71.1	Min. : 52.0
1st Qu.:120.8	1st Qu.: 96.5
Median :196.3	Median :123.0
Mean :230.7	Mean :146.7
3rd Qu.:326.0	3rd Qu.:180.0
Max. :472.0	Max. :335.0

qsec	vs
Min. :14.50	Min. :0.0000
1st Qu.:16.89	1st Qu.:0.0000
Median :17.71	Median :0.0000
Mean :17.85	Mean :0.4375
3rd Qu.:18.90	3rd Qu.:1.0000
Max. :22.90	Max. :1.0000

carb
Min. :1.000
1st Qu.:2.000
Median :2.000
Mean :2.812
3rd Qu.:4.000
Max. :8.000

mpg	Miles/(US) gallon
cyl	Number of cylinders
disp	Displacement (cu.in.)
hp	Gross horsepower
drat	Rear axle ratio
wt	Weight (1000 lbs)
qsec	1/4 mile time
vs	Engine (0 = V-shaped, 1 = straight)
am	Transmission (0 = automatic, 1 = manual)
gear	Number of forward gears
carb	Number of carburetors

Removed (binary/categoric data)

## When they are not centered and scaled, fewer PC dominate and mask the effect of other PCs

```
> mtcarsCenteredScaled.pca <- prcomp(mtcars[,c(1:7,10,11)], center = TRUE, scale. = TRUE)
> summary(mtcarsCenteredScaled)
Error in summary(mtcarsCenteredScaled) :
  object 'mtcarsCenteredScaled' not found
> summary(mtcarsCenteredScaled.pca)
Importance of components:
```

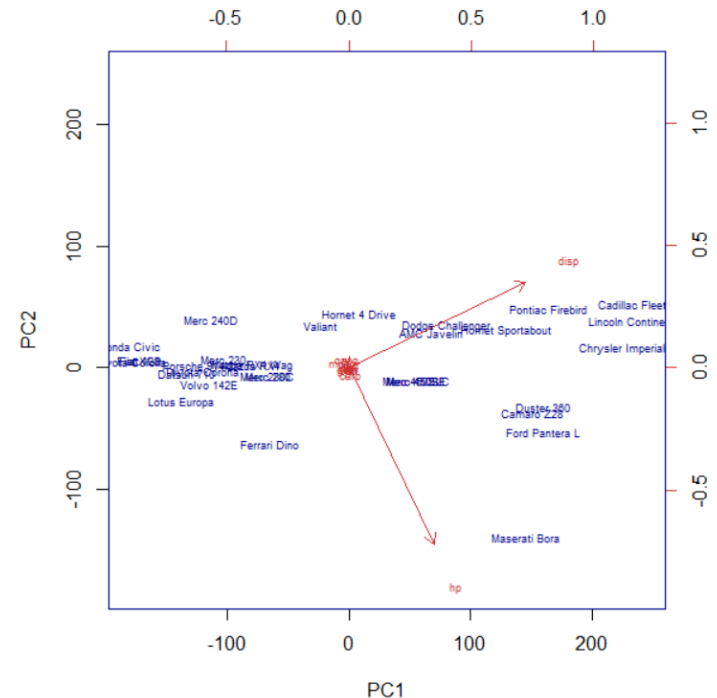
	PC1	PC2	PC3	PC4	PC5	PC6	PC7	PC8	PC9
Standard deviation	2.3782	1.4429	0.71008	0.51481	0.42797	0.35184	0.32413	0.2419	0.14896
Proportion of Variance	0.6284	0.2313	0.05602	0.02945	0.02035	0.01375	0.01167	0.0065	0.00247
Cumulative Proportion	0.6284	0.8598	0.91581	0.94525	0.96560	0.97936	0.99103	0.9975	1.00000

```
> mtcarsCenteredNOTScaled.pca <- prcomp(mtcars[,c(1:7,10,11)], center = TRUE, scale. = FALSE)
> summary(mtcarsCenteredNOTScaled.pca)
Importance of components:
```

	PC1	PC2	PC3	PC4
Standard deviation	136.532	38.14735	3.06642	1.27
Proportion of Variance	0.927	0.07237	0.00047	0.00000
Cumulative Proportion	0.927	0.99938	0.99985	0.99985

```
> summary(mtcarsNOTCenteredNOTScaled.pca)
Importance of components:
```

	PC1	PC2	PC3
Standard deviation	310.1169	40.88309	15.83774
Proportion of Variance	0.9803	0.01704	0.00256
Cumulative Proportion	0.9803	0.99738	0.99993



## Cpptraj can do it?

We couldn't find a way of transforming the coordinates like this:

$$x'_i = \frac{(x_i - \bar{x})}{std.dev_i}$$

Nor could find a way of calling diagmatrix in cpptraj so that it does it automatically.  
Neither we could find a way to do it using GROMACS



## Our approach: DIY

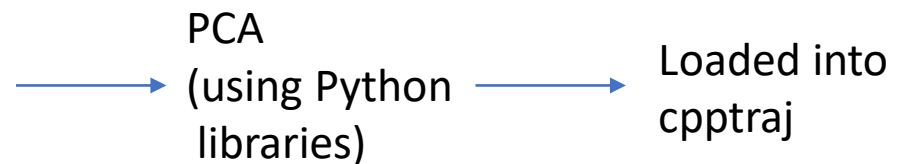
We did export each “frame” of the simulation in text-readable format:

- 2500 PDB files (23 Mb), each with two proteins and many chlorine, sodium, lipids, Water molecules...

We removed chlorine, sodium, lipids, water, the protein “B”, and only kept Carbon-alpha of protein A.

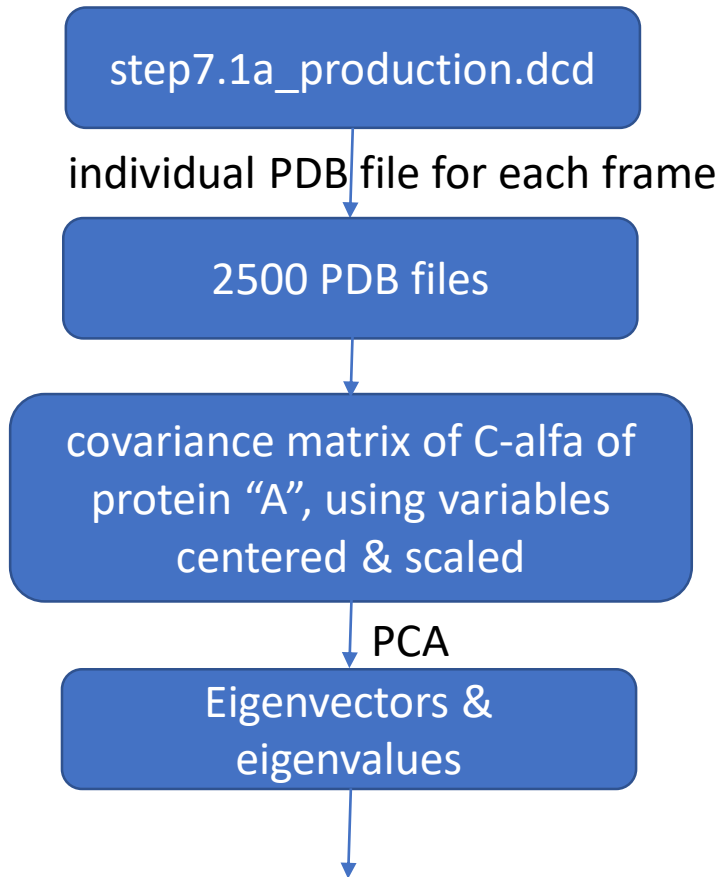
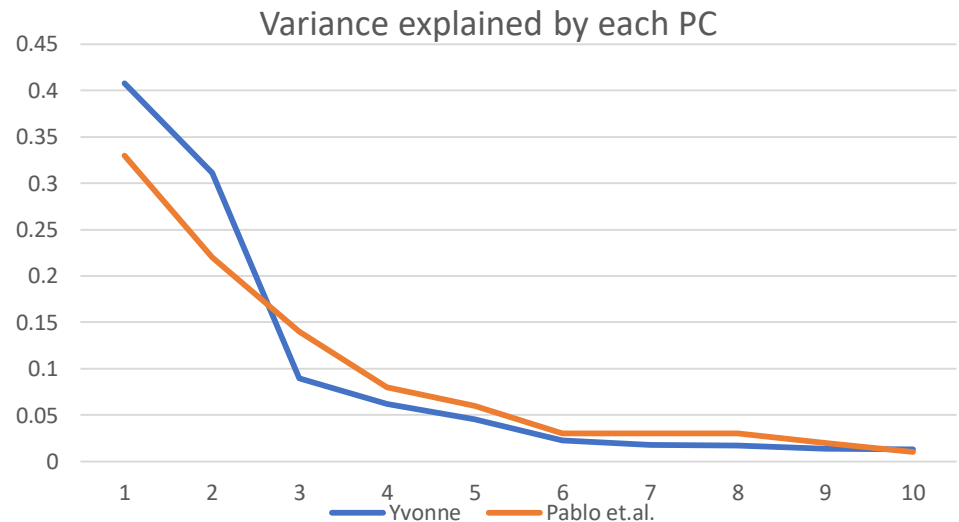
Prepared a table where each cell's value was centered & scaled

	x1	y1	z1	x2	y2	z2
T1						
T2						
...						
T2500						



## Our approach: DIY

With centering & scaling,  
more principal components are needed  
to describe the C-alfa movements in the  
simulation (more information is captured by PCs)

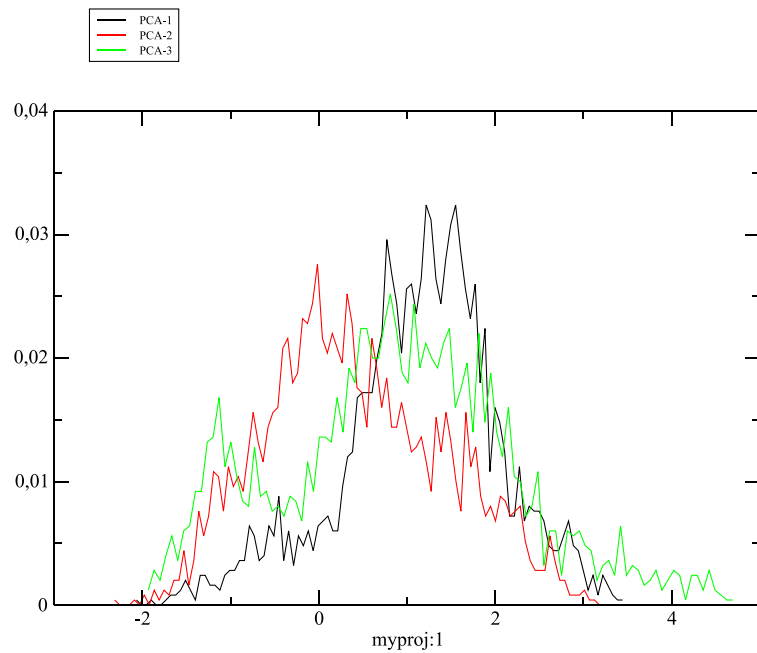


Convert eigenvector and eigenvalues  
into cpptraj format and load them into  
Cpptraj, and do the same analysis after PCA



# Our results: histogram of projection using PCA eigenvectors

## With our scaling



## Original (without scaling)

