

# Real dataset tutorial – PCA plot

## Goal of this real dataset tutorial

In this tutorial using real dataset, we want to see if we can answer the following biological question: Can we see the pattern of clustering based on different proportions of barley grain fed in the cows?

## Dataset description

“Metabolite concentrations of 39 rumen samples measured by proton NMR from dairy cows fed with different proportions of barley grain (Ametaj BN, et al.). Group label - 0, 15, 30, or 45 - indicating the percentage of grain in diet.”

## Step by step installation and running

Let's get started! First, we need to install the package.

Steps: 1. Git clone or download the github folder;

2. Open the terminal, and go to this folder;

3. Enter

`pip install dist/metabolomics_analysis_tools-0.1.0.tar.gz` to install the package locally;

Then, we can import functions we will use for this demo from the package `metabolomics_analysis_tools` @import

```
import metabolomics_analysis_tools.data_preprocessing.data_reading as dr
import metabolomics_analysis_tools.data_preprocessing.normalization as dn
import metabolomics_analysis_tools.stats_analyses.analyses as sa
import warnings
warnings.filterwarnings('ignore')
```

1. Then we can use the `data_reading` module to read in the data, by default it will read in the data from the `resources/test_dataset` folder in the package.

We can also use the `data_reading` module to read in the data from a custom path, by passing the path as an argument to the `read_data_file` function (`file_path='path/to/file.csv'`).

The `read_data_file` function will return a pandas dataframe. We would also want to get rid off 0 value in the data to avoid division by zero issue.

```
df=dr.read_data_file(file_path='resources/test_dataset/cow_diet.csv')
df = df.replace(0, 0.0001)
df.head()
```

data read successfully  
the shape of the dataframe is: (39, 49)

	Sample	Diet	1,3-D	3-HB	3-HP	3-PP	Acetate	Acetoacetate	Alanine	Aspartate	Benzoa
0	0_1_1	0	2.5	0.0001	11.2	312.5	43720.6	12.8	91.3	62.9	23
1	0_1_3	0	2.3	47.7000	35.6	465.5	55083.7	51.8	153.0	133.9	24
2	0_1_5	0	3.3	34.2000	18.8	363.4	39246.2	20.2	86.5	95.1	11
3	0_1_7	0	4.1	29.9000	8.5	316.5	40485.1	25.2	107.3	55.8	14
4	0_1_10	0	10.9	42.5000	50.8	680.4	53203.9	55.1	247.9	178.9	63

- Next we can use the normalization module to normalize the data, here we will use the median normalization method `normalized_data=dn.normalize_by_median(df)`.

We can have a look at the first 5 rows of the normalized data `normalized_data.head()`.

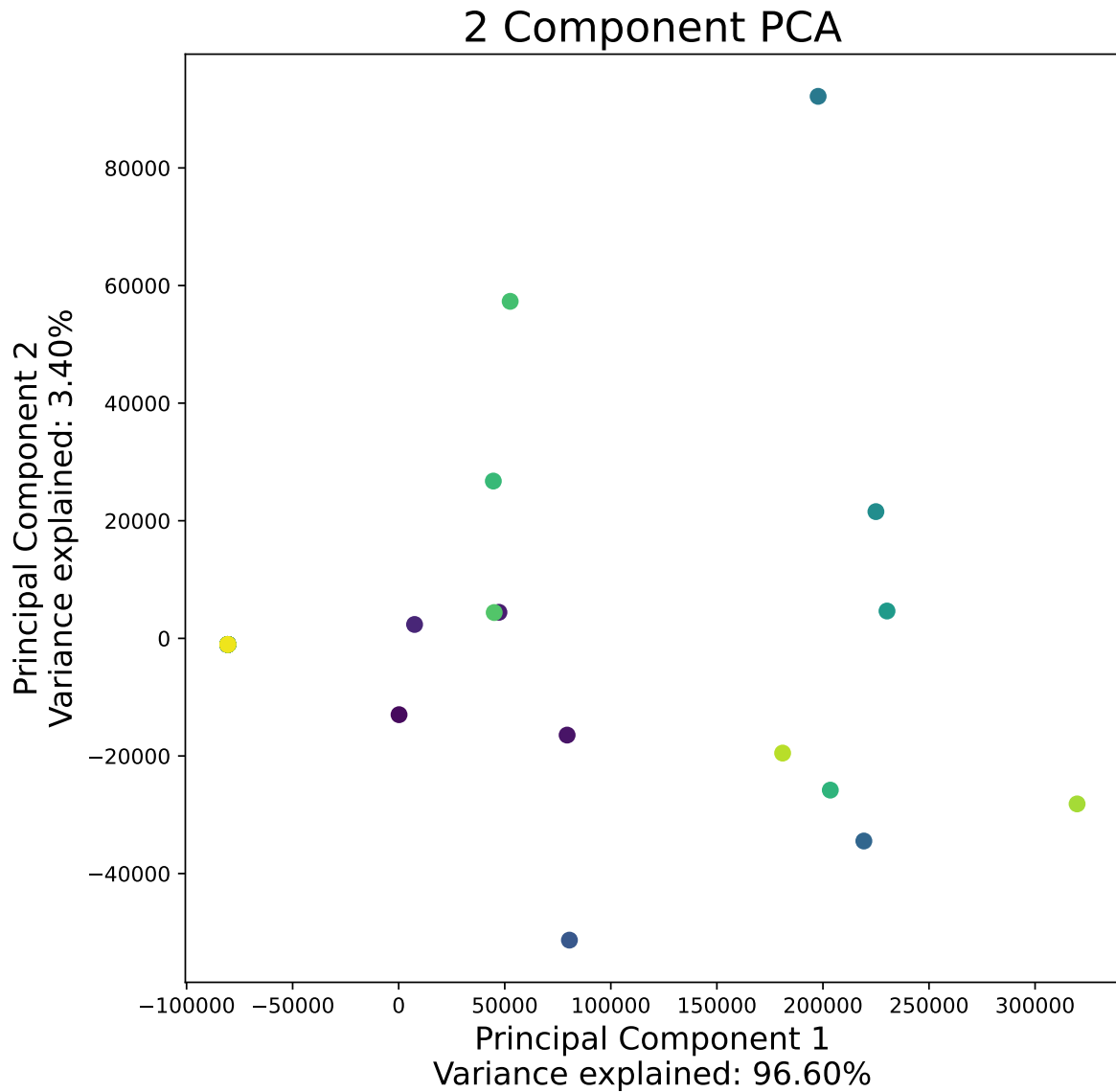
```
normalized_data=dn.normalize_by_median(df)
normalized_data.head()
```

	Sample	Diet	1,3-D	3-HB	3-HP	3-PP	Acetate	Acetoacetate	Alanine	Aspart
0	0_1_1	0	0.416667	0.000002	0.273171	0.935349	0.927099	0.231047	0.283101	0.469
1	0_1_3	0	0.383333	0.979466	0.868293	1.393295	1.168055	0.935018	0.474419	1.000
2	0_1_5	0	0.550000	0.702259	0.458537	1.087698	0.832219	0.364621	0.268217	0.710
3	0_1_7	0	0.683333	0.613963	0.207317	0.947321	0.858490	0.454874	0.332713	0.416
4	0_1_10	0	1.816667	0.872690	1.239024	2.036516	1.128193	0.994585	0.768682	1.336

- We can use the analyses module to perform statistical analyses on the data.  
Here we will first perform a PCA analysis on the data to see if there are any patterns in the data.

The `PCA_analysis` function will return a pandas dataframe containing the principal components `principal_components=sa.PCA_analysis(normalized_data)`.

```
principal_components=sa.PCA_analysis(normalized_data)
```



- Interpretation of the PCA results:

From the PCA results, we can see it is hard to distinguish the groups of cows based on their different feeding plan (indicated by colors). However, the first two PCs explained most of the variance in the data. Therefore, even though the clustering didn't give us a good pattern for different groups of data, I might assume that it was due to limited number of samples. If we have more samples, we might be able to see a better pattern.