

Multivariate_Analysis_Metabolomics_Workflow

September 22, 2022

1 Introduction

The purpose of this notebook was to illustrate the data analysis performed along my MSc thesis, which included unsupervised as well as supervised Multivariate Analysis (MVA) carried out on a metabolomic dataset. It was also shown how tools namely Probabilistic Quotient Normalization (PQN) or Variable Importance in Projection (VIP) were implemented in Python. Lastly, it was shown how to analyse MetaboRank output and perform feature mass-matching to find hits.

The dataset used in this project came from a neurotoxicology experiment where astrocyte cells were exposed to different concentration of digoxin. Digoxin is one of the oldest cardiovascular medication used nowadays; digoxin is a common agent used to manage atrial fibrillation and the symptoms of heart failure as it is a positive inotropic and negative chronotropic drug which means that digoxin increases the force of the heartbeat and decreases the heart rate.

The digoxin dataset corresponded to 6 experimental groups by 4 replicates, leading to a total of twenty-four samples. The experimental groups could be summarized as *Negative Control*, *Exposed Groups*, and *Positive Control*. On the Exposed groups the concentration of digoxin ranged from 0 to 10 μmolar . The Positive Control group was exposed to $\text{TNF}\alpha$.

The different steps of the data analysis carried out in this project can be summarized as follow:

- Exploratory Data Analysis
 - Principal Component Analysis (PCA)
 - Hierarchical Clustering Analysis (HCA)
- Supervised modelling
 - Partial Least Squares (PLS)
 - Random Forest (RF)
- Metabolic Fingerprint Extraction
- MetaboRank Output Analysis
- Metabolite - Feature Mass Matching

Author: [Christian Peralta](#)

1.1 Importing libraries

```
[1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import plotly.express as px
import plotly.graph_objects as go
```

```
from plotly.subplots import make_subplots
from xlrd import open_workbook_xls
```

2 Data Preparation

The original dataset contained several column which were not of interest for a first exploratory data analysis (EDA) where regardless of the model to use, the best practice is to have a data matrix containing only [observations x predictors]. In addition to perform data cleaning, observations label was modified to a shorter but meaningful name, enabling a straightforward interpretation of the results.

```
[2]: df = pd.read_csv("20201209_MGZ_720_RP.csv", sep= "\t", header= 4)
df.shape
```

```
[2]: (3255, 71)
```

```
[3]: df.columns
```

```
[3]: Index(['Compound', 'Neutral mass (Da)', 'm/z', 'Charge',
          'Retention time (min)', 'Chromatographic peak width (min)',
          'Identifications', 'Anova (p)', 'q Value', 'Max Fold Change',
          'Highest Mean', 'Lowest Mean', 'Isotope Distribution',
          'Maximum Abundance', 'Minimum CV%', 'MSMS info available', 'Rt',
          'KEGG in silico', 'Waters CCS library', 'LipidBlast', 'A revoir',
          'Antibiotic', 'pH indicator', 'Accepted ID', 'Accepted Compound ID',
          'Accepted Description', 'Adducts', 'Formula', 'Score',
          'Fragmentation Score', 'Mass Error (ppm)', 'Isotope Similarity',
          'Retention Time Error (mins)', 'Compound Link', 'Privileged',
          'Control A_1-A,1_01_201', 'Control B_1-A,2_01_210',
          'Control C_1-A,3_01_216', 'Control D_1-A,4_01_191',
          'Digoxine conc 0.0 A _1-B,1_01_214',
          'Digoxine conc 0.0 B _1-B,2_01_218', 'Digoxine conc 0.0 C_1-B,3_01_190',
          'Digoxine conc 0.0 D _1-B,4_01_202', 'Digoxine conc 0.1 A_1-C,1_01_205',
          'Digoxine conc 0.1 B_1-C,2_01_209', 'Digoxine conc 0.1 C_1-C,3_01_198',
          'Digoxine conc 0.1 D_1-C,4_01_192', 'Digoxine conc 1.0 A_1-D,1_01_208',
          'Digoxine conc 1.0 B_1-D,2_01_200', 'Digoxine conc 1.0 C_1-D,3_01_193',
          'Digoxine conc 1.0 D_1-D,4_01_215', 'Digoxine conc 10.0 A_1-E,1_01_213',
          'Digoxine conc 10.0 B_1-E,2_01_206',
          'Digoxine conc 10.0 C_1-E,3_01_199',
          'Digoxine conc 10.0 D _1-E,4_01_194', 'TNF-alfa A_1-F,1_01_207',
          'TNF-alfa B_1-F,2_01_217', 'TNF-alfa C_1-F,3_01_197',
          'TNF-alfa D_1-F,4_01_189', 'dQC_1-A,8_01_186', 'dQC_1-A,8_01_188',
          'dQC_1-A,8_01_196', 'dQC_1-A,8_01_204', 'dQC_1-A,8_01_212',
          'QC_1-A,7_01_185', 'QC_1-A,7_01_187', 'QC_1-A,7_01_195',
          'QC_1-A,7_01_203', 'QC_1-A,7_01_211', 'QC_1-A,7_01_219',
          'QC_1-A,7_01_221'],
          dtype='object')
```

The data frame was cleaned to select those column corresponding to the experimental samples.

```
[4]: df2 = df.iloc[:, 35:]
df2.columns
```

```
[4]: Index(['Control A_1-A,1_01_201', 'Control B_1-A,2_01_210',
          'Control C_1-A,3_01_216', 'Control D_1-A,4_01_191',
          'Digoxine conc 0.0 A _1-B,1_01_214',
          'Digoxine conc 0.0 B _1-B,2_01_218', 'Digoxine conc 0.0 C_1-B,3_01_190',
          'Digoxine conc 0.0 D _1-B,4_01_202', 'Digoxine conc 0.1 A_1-C,1_01_205',
          'Digoxine conc 0.1 B_1-C,2_01_209', 'Digoxine conc 0.1 C_1-C,3_01_198',
          'Digoxine conc 0.1 D_1-C,4_01_192', 'Digoxine conc 1.0 A_1-D,1_01_208',
          'Digoxine conc 1.0 B_1-D,2_01_200', 'Digoxine conc 1.0 C_1-D,3_01_193',
          'Digoxine conc 1.0 D_1-D,4_01_215', 'Digoxine conc 10.0 A_1-E,1_01_213',
          'Digoxine conc 10.0 B_1-E,2_01_206',
          'Digoxine conc 10.0 C_1-E,3_01_199',
          'Digoxine conc 10.0 D _1-E,4_01_194', 'TNF-alfa A_1-F,1_01_207',
          'TNF-alfa B_1-F,2_01_217', 'TNF-alfa C_1-F,3_01_197',
          'TNF-alfa D_1-F,4_01_189', 'dQC_1-A,8_01_186', 'dQC_1-A,8_01_188',
          'dQC_1-A,8_01_196', 'dQC_1-A,8_01_204', 'dQC_1-A,8_01_212',
          'QC_1-A,7_01_185', 'QC_1-A,7_01_187', 'QC_1-A,7_01_195',
          'QC_1-A,7_01_203', 'QC_1-A,7_01_211', 'QC_1-A,7_01_219',
          'QC_1-A,7_01_221'],
          dtype='object')
```

As one can see, above, the names were long. The name was given following a certain experimental regulation which is not of interest for data analysis.

```
[5]: new_names = ['Control_A', 'Control_B', 'Control_C', 'Control_D', 'Conc_0.0_A',
                  ↪ 'Conc_0.0_B', 'Conc_0.0_C', 'Conc_0.0_D',
                  'Conc_0.1_A', 'Conc_0.1_B', 'Conc_0.1_C', 'Conc_0.1_D', 'Conc_1.
                  ↪ 0_A', 'Conc_1.0_B', 'Conc_1.0_C', 'Conc_1.0_D',
                  'Conc_10.0_A', 'Conc_10.0_B', 'Conc_10.0_C', 'Conc_10.0_D',
                  ↪ 'TNF-a_A', 'TNF-a_B', 'TNF-a_C', 'TNF-a_D', 'dQC_1',
                  'dQC_2', 'dQC_3', 'dQC_4', 'dQC_5', 'QC_1', 'QC_2', 'QC_3', 'QC_4',
                  ↪ 'QC_5', 'QC_6', 'QC_7']
```

```
[6]: # using set_axis for a clean axis manipulation
df2 = df2.set_axis(new_names, axis=1)
```

Once the dataset was filtered the last step to be ready for EDA was to transpose the dataset, hence it was switched from [predictors x observations] to [observations x predictors]. Where the observations were stored as the index of the pandas.DataFrame instance.

```
[7]: # transposing the data frame
df3 = df2.T

# defining column name for the index
```

```
df3.index.name = "samples"

# checking the dimensions
df3.shape
```

[7]: (36, 3255)

3 Exploratory Data Analysis

Typical multivariate analysis workflow in metabolomics starts with exploratory data analysis (EDA) which consists on unsupervised modeling to provide a first overview of the data. Unsupervised statistical tools aim at building models summarizing the dataset in an intelligible manner and hopefully finding natural partitions of the dataset to facilitate the understanding of the relationship between the samples and detect potential outliers. The models also provide information about the variables that are responsible for these relationships.

The EDA can be seen as an iterative process where sample groups were excluded from the model at each iteration until keeping only exposure-related samples whilst understanding the data. The first PCA model was built using the entire dataset, which included QCs and dQCs, aiming at having a first overview of the data focused on studying the quality of the data acquisition.

3.1 How to perform Principal Component Analysis in Python

```
[8]: from sklearn.decomposition import PCA
     from sklearn.preprocessing import StandardScaler
```

The data must be standardize before PCA

```
[9]: #normalizing the data
     df_pca = StandardScaler().fit_transform(df3)
```

```
[10]: pca = PCA(n_components=.95)
      Pcomponents = pca.fit_transform(df_pca)
```

```
[11]: # computing explaining variable
     var = pca.explained_variance_ratio_*100
```

The scores are returned directly from sklearn function, however, loadings have to be manually computed from sklearn output.

To obtain loadings, the loading matrix (also known as cross-correlation matrix) which is given by the equation $loadings = eigenvectors * \sqrt{eigenvalues}$ and can be done with sklearn as follows:

```
loadings = pca.components_.T * np.sqrt(pca.explained_variance_)
```

```
[12]: loadings = pd.DataFrame(pca.components_.T * np.sqrt(pca.explained_variance_))
```

3.2 Complete dataset

It was defined four list of names to help with visualization and interpretation of the results by using them to set color and shape of markers as well as color and symbol sequence

```
[13]: color = ["control", "control", "control", "control", "[C] 0.0", "[C] 0.0", "[C] 0.0", "[C] 0.0", "[C] 0.1", "[C] 0.1", "[C] 0.1", "[C] 0.1", "[C] 1.0", "[C] 1.0", "[C] 1.0", "[C] 1.0", "[C] 10.0", "[C] 10.0", "[C] 10.0", "[C] 10.0", "TNF-a", "TNF-a", "TNF-a", "dQC", "dQC", "dQC", "dQC", "dQC", "dQC", "QC", "QC", "QC", "QC", "QC", "QC", "QC"]

color_seq=["mediumblue", "darkorange", "limegreen", "orchid", "red", "maroon", "navy", "grey"]

symbol = ["negative group", "negative group", "negative group", "negative group", "no-exposure", "no-exposure", "no-exposure", "no-exposure", "low-exposure", "low-exposure", "low-exposure", "low-exposure", "mid-exposure", "mid-exposure", "mid-exposure", "mid-exposure", "high-exposure", "high-exposure", "high-exposure", "high-exposure", "positive group", "positive group", "positive group", "positive group", "QC group", "QC group", "QC group", "QC group", "QC group", "QC group", "QC group", "QC group"]

symbol_seq=["circle", "diamond", "triangle-up", "triangle-down", "x", "diamond-wide", "pentagon"]
```

Score plot visualization

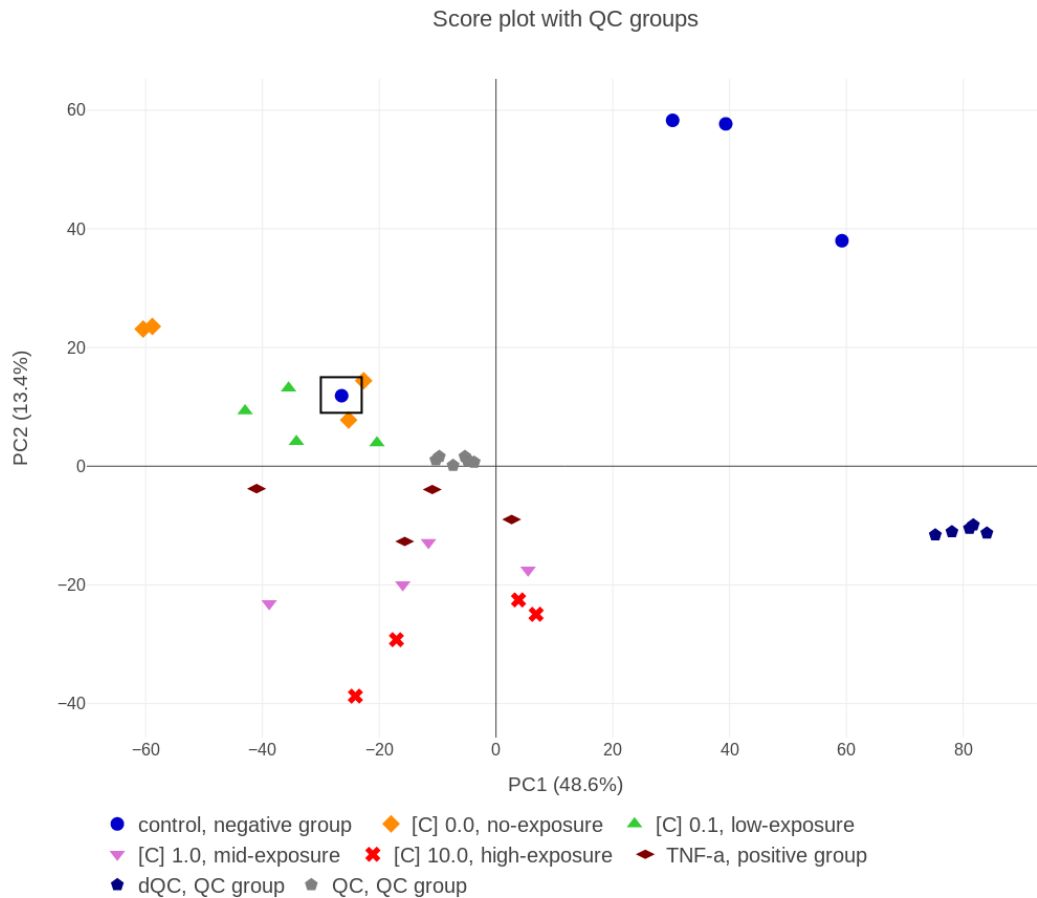
```
[14]: score_plot = px.scatter(Pcomponents, x= Pcomponents[:,0], y=Pcomponents[:,1], labels = {"x" : f"PC1 ({var[0]:.1f}%)", "y" : f"PC2 ({var[1]:.1f}%)", color=color, template= "presentation", width= 900, height= 900, title= "Score plot with QC groups", hover_name= df3.index, symbol=symbol, symbol_sequence=symbol_seq, color_discrete_sequence=color_seq)

score_plot.update_layout(legend_title=None, font=dict(size=16), legend=dict(yanchor="top", y=-.1, xanchor="left", x=0.01, orientation="h", font_size=20))

score_plot.update_traces(marker=dict(size=13))

score_plot.add_shape(type="rect", x0=-30, x1=-23, y0=9, y1=15)

score_plot.show()
```

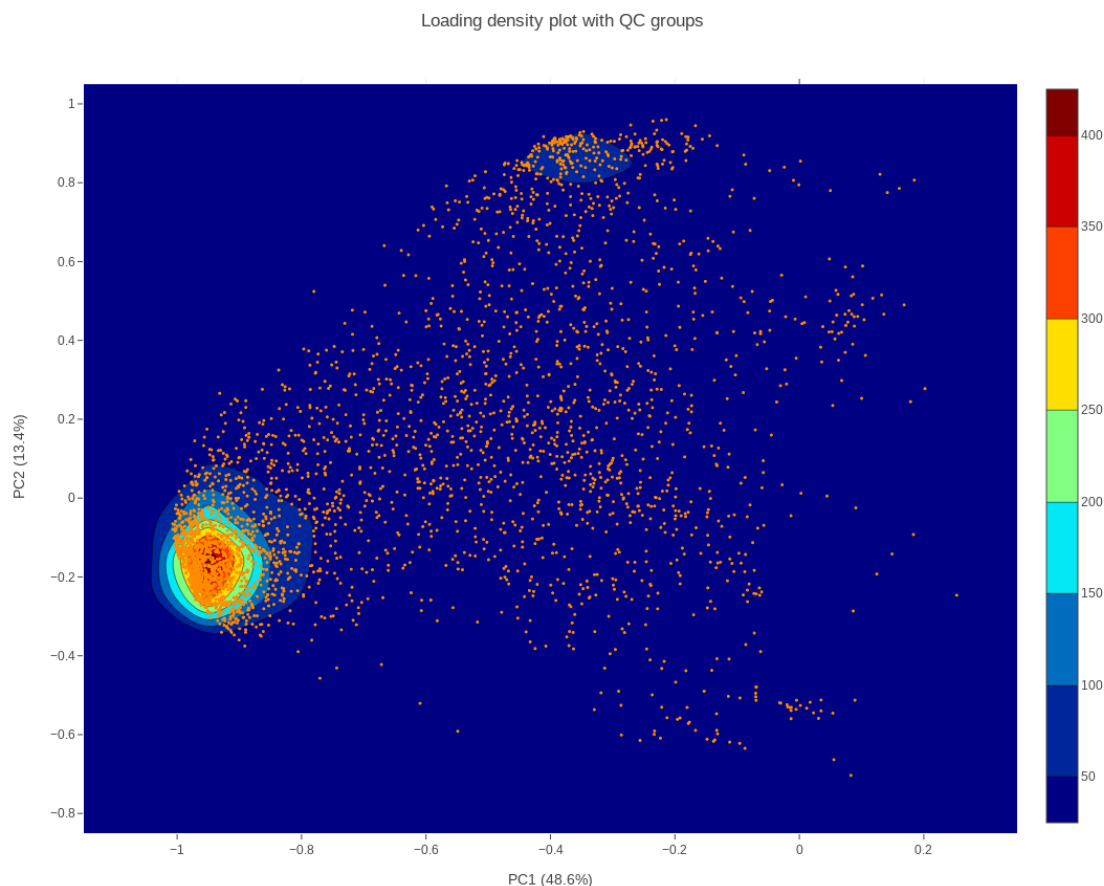


PCA model was built on the complete dataset where the above score plot, showed how PC1 was driven by the dQCs samples, whereas PC2 was mainly driven by the control group, which seemed to be orthogonal to digoxin effect. In addition, on the PC2 the trending of digoxin concentration was reflected; on top there were the samples with none or lower digoxin concentration, followed by samples with the higher concentration. One sample belonging to the control group - highlighted within a rectangle - was found in between [C] 0.0 and [C] 0.1 samples, hence it was considered as an outlier. In the middle of the trending, there were located the positive control group, e.g. TNF- , and close to the origin one could find QCs samples which were a pooled mixture of all the samples and acted as a biological mean.

Loading plot visualization

```
[15]: loading_plot = go.Figure()
loading_plot.add_trace(go.Scatter(x=loadings.iloc[:,0], y=loadings.iloc[:,1],
    ↪ mode="markers", marker=dict(size=3, color="darkorange"), xaxis="x",
    ↪ yaxis="y"))
loading_plot.add_histogram2dcontour(x=loadings.iloc[:,0], y=loadings.iloc[:,1],
    ↪ colorscale="Jet", xaxis="x", yaxis="y")
```

```
loading_plot.update_layout(template="presentation", width= 900, height= 900,
    ↪hovermode=False, font=dict(size=12), title="Loading density plot with QC_
    ↪groups")
loading_plot.update_xaxes(title_text=f"PC1 ({var[0]:.1f}%)")
loading_plot.update_yaxes(title_text=f"PC2 ({var[1]:.1f}%)")
```



Through the loadings plot, it was confirmed that the PC1 was driven by dQC, as aforementioned, since the concentration of the loading at the positive side was notably lower than at the negative side. This was related to the fact that dQC were a diluted mixture of all samples, hence containing lower concentration of compounds.

After a first exploration of the data, it was highlighted that the data was acquired with high quality as the QCs and dQCs samples were well grouped. Moreover, one negative control sample was considered as an outlier and that the negative control group seemed to behave orthogonally to digoxin effect. Consequently, next exploratory steps included removal of QCs/dQCs to observe how the first 2 PCs behaved.

3.3 dataset excluding: *QCs/dQCs*

```
[147]: # using pd.filter method with negative regex matching to filter out QC samples
df_filtered = df3.filter(axis=0, regex="(?!m)(?!dQC|QC).)")

[148]: color = ["control", "control", "control", "control", "[C] 0.0", "[C] 0.0", "[C] 0.0", "[C] 0.0", "[C] 0.1", "[C] 0.1",
               "[C] 0.1", "[C] 0.1", "[C] 1.0", "[C] 1.0", "[C] 1.0", "[C] 1.0", "[C] 10.0", "[C] 10.0", "[C] 10.0", "[C] 10.0",
               "[C] 10.0", "[C] 10.0", "TNF-a", "TNF-a", "TNF-a", "TNF-a"]

color_seq=["mediumblue", "darkorange", "limegreen", "orchid", "red", "maroon"]

symbol = ["negative group", "negative group", "negative group", "negative group", "no-exposure", "no-exposure", "no-exposure", "no-exposure",
          "low-exposure", "low-exposure", "low-exposure", "low-exposure", "mid-exposure", "mid-exposure", "mid-exposure", "mid-exposure",
          "high-exposure", "high-exposure", "high-exposure", "high-exposure", "positive group", "positive group", "positive group", "positive group"]

symbol_seq=["circle", "diamond", "triangle-up", "triangle-down", "x", "diamond-wide", "pentagon"]

[149]: #Let's create a pipeline to concatenate scaler and pca
from sklearn.pipeline import Pipeline

pca_out= PCA(n_components=.95)
pca_pipe = Pipeline([("scale", StandardScaler()), ("pca", pca_out)])

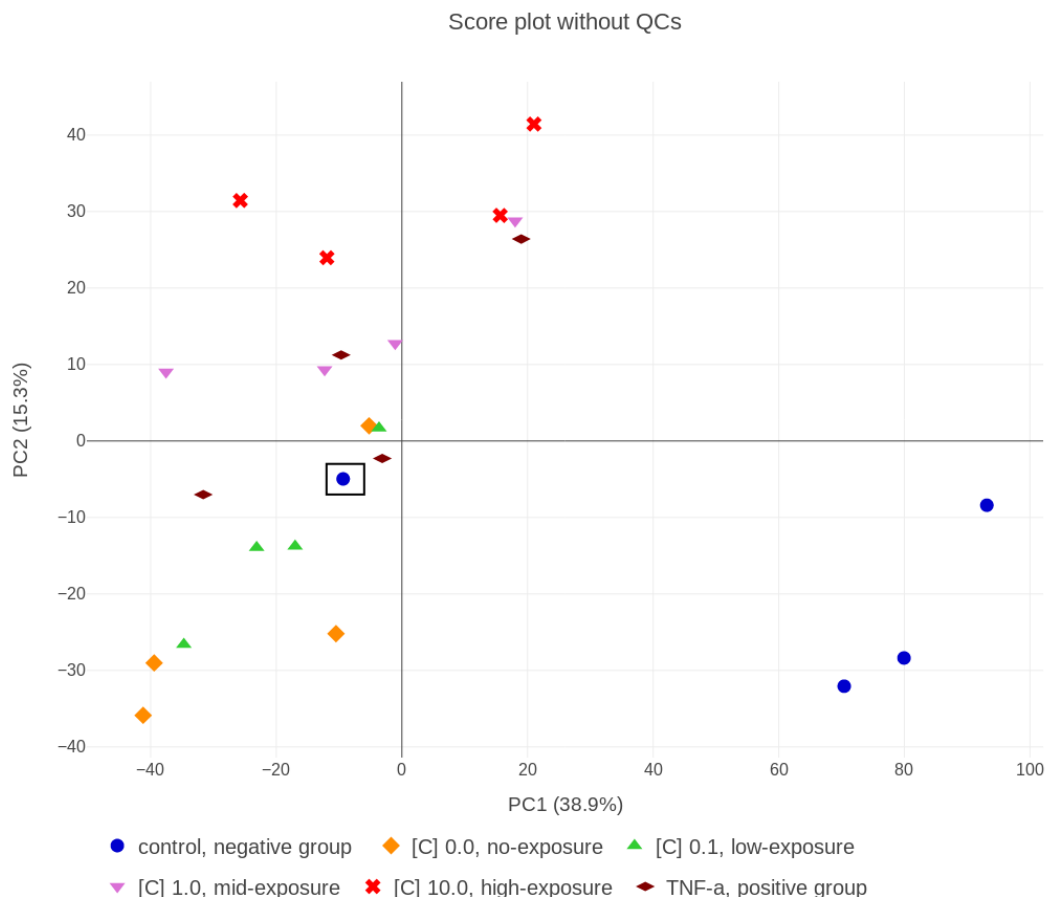
[150]: pcomponents = pca_pipe.fit_transform(df_filtered)
var = pca_out.explained_variance_ratio_*100
```

Score plot visualization

```
[151]: score_plot = px.scatter(data_frame= pcomponents, x=pcomponents[:,0], y=pcomponents[:,1], color=color,
                               labels= {"x": f"PC1 ({var[0]:.1f}%)", "y": f"PC2 ({var[1]:.1f}%)"}, template= "presentation", width= 900, height= 900,
                               title= "Score plot without QCs", hover_name= df_filtered.index, symbol=symbol, symbol_sequence=symbol_seq, color_discrete_sequence=color_seq)
score_plot.update_layout(legend_title= None, font=dict(size=16), legend=dict(yanchor="top", y=-.1, xanchor="left", x=0.01, orientation="h", font_size=20), overwrite=True)
score_plot.update_traces(marker=dict(size=13))
score_plot.add_shape(type="rect", x0=-12, x1=-6, y0=-7, y1=-3)
```



```
score_plot.show()
```

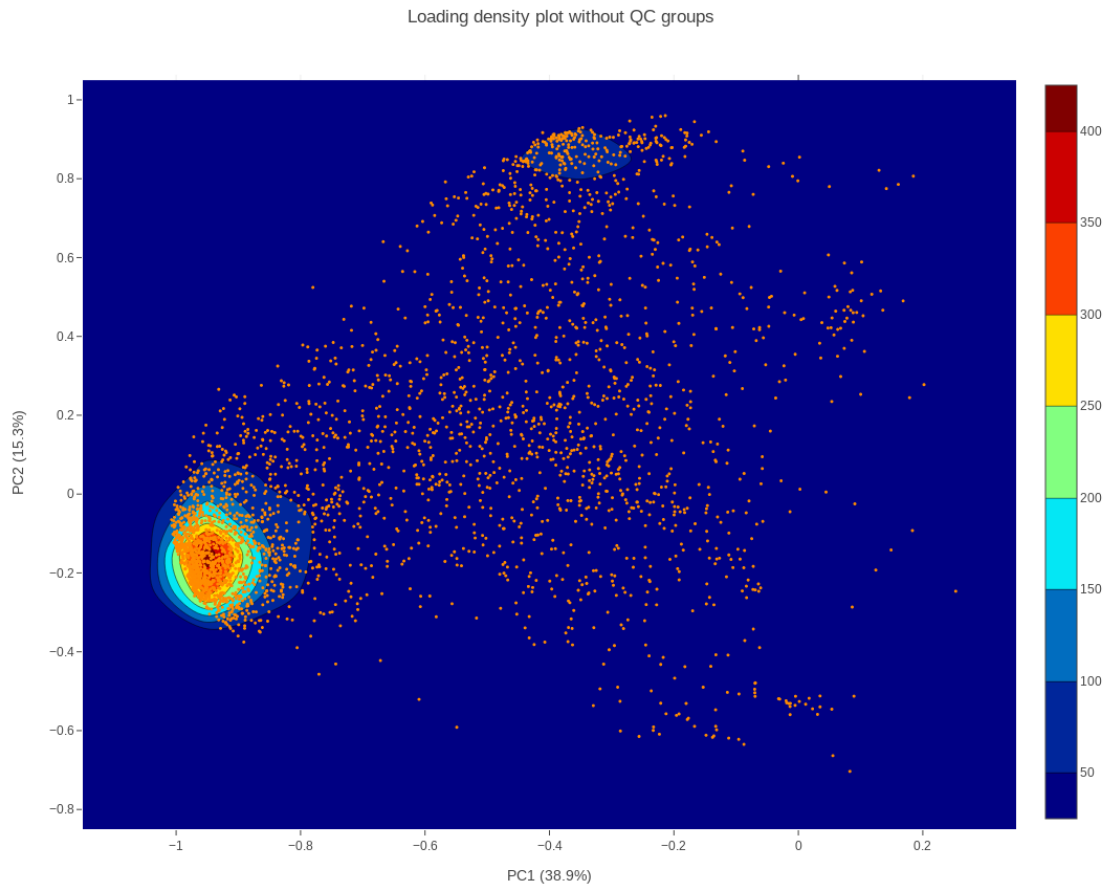


Exclusion of QC groups from the dataset, allowed a better visualization of the digoxin trending by the score plot. At the same time, the orthogonality of the negative group to the toxicity of digoxin was plain to see. However, it was hard to interpret what was driving the PCs as the trending in concentration was crossing both PCs rather than been represented by only one.

Loading plot visualization

```
[21]: loading_plot = go.Figure()
loading_plot.add_trace(go.Scatter(x=loadings.iloc[:,0], y=loadings.iloc[:,1],
    ↪ mode="markers", marker=dict(size=3, color="darkorange"), xaxis="x",
    ↪ yaxis="y"))
loading_plot.add_histogram2dcontour(x=loadings.iloc[:,0], y=loadings.iloc[:,1],
    ↪ colorscale="Jet", xaxis="x", yaxis="y")
loading_plot.update_layout(template="presentation", width= 900, height= 900,
    ↪ hovermode=False, font=dict(size=12), title="Loading density plot without QC_
    ↪ groups")
```

```
loading_plot.update_xaxes(title_text=f"PC1 ({var[0]:.1f}%)")
loading_plot.update_yaxes(title_text=f"PC2 ({var[1]:.1f}%)")
loading_plot.show()
```



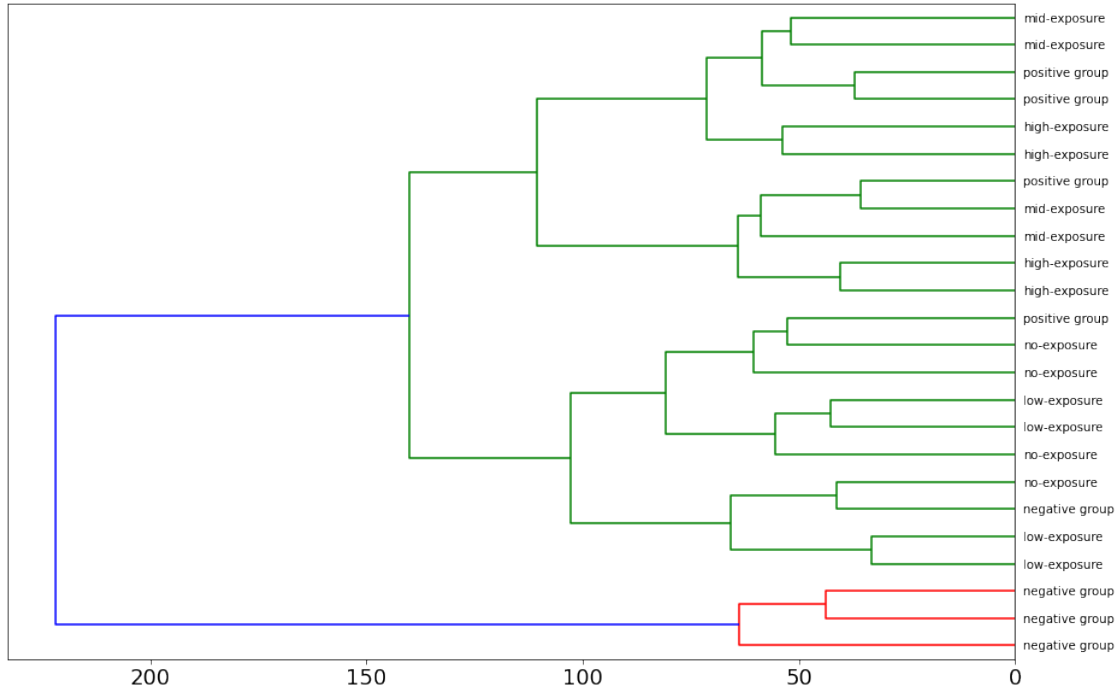
On the other hand, the loading plot, Figure 3.3b, showed an unusual behavior of the loadings. It was observed a hotspot located on the left side of the PC1 which was not related to any sample or group of samples in particular as it was the case in the first PCA model.

HCA was performed on the PCA components to prove what was observed, that control group was orthogonal to the rest of the groups.

```
[22]: import scipy.cluster.hierarchy as sch
```

```
[152]: plt.figure(figsize=(15,10))

        dendrogram_pca = sch.dendrogram(sch.linkage(pcomponents, method='ward'),
        ↪ labels=symbol, orientation="left")
        plt.show()
```



Hence, to continue with the EDA, the negative control group was removed from subsequent PCA models.

3.4 Dataset excluding: *QCs/dQCs, negative control groups*

```
[24]: df_filtered = df_filtered.filter(axis=0,regex="(?!^Control*)")

[25]: color = ["[C] 0.0", "[C] 0.0", "[C] 0.0", "[C] 0.0", "[C] 0.1", "[C] 0.1",
              "[C] 0.1", "[C] 0.1", "[C] 1.0", "[C] 1.0", "[C] 1.0", "[C] 1.0",
              "[C] 10.0", "[C] 10.0", "[C] 10.0", "[C] 10.0", "TNF-a", "TNF-a", "TNF-a", "TNF-a"]

color_seq=["darkorange", "limegreen", "orchid", "red", "maroon"]

symbol = ["no-exposure", "no-exposure", "no-exposure", "no-exposure",
          "low-exposure", "low-exposure", "low-exposure", "low-exposure",
          "mid-exposure", "mid-exposure", "mid-exposure", "mid-exposure",
          "high-exposure", "high-exposure", "high-exposure", "high-exposure",
          "positive group", "positive group", "positive group", "positive group"]

symbol_seq=["diamond", "triangle-up", "triangle-down", "x", "diamond-wide", "pentagon"]
```

From now on, the analysis was done using a Python library developed during this project to speed up data analysis workflow. Hence, following steps also show how to use `chemometrics` library.

```
[26]: from chemometrics.unsupervised import unsupervised
```

```
[27]: # creating the object
      output = unsupervised()

      # computing pca, pc variable contains the scores
      pca = output.PCA_ready(df = df_filtered)
```

Score plot visualization

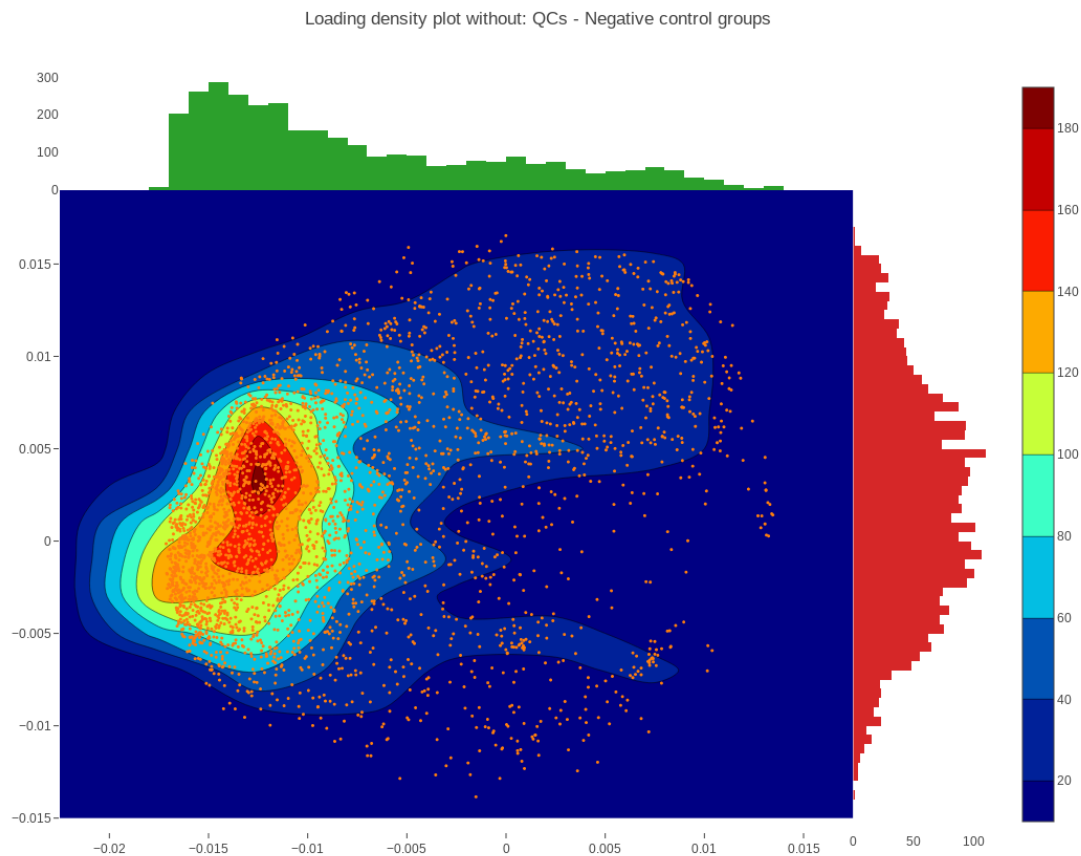
```
[28]: # now the object containin the results of pca has to be used, here named
      ↪ 'output'
      score_plot = output.score_viz(pcx= 0, pcy= 1, color=color, symbol=symbol,
      ↪ label= None, label_position = None, symbol_sequence=symbol_seq,
      color_discrete_sequence=color_seq,
      ↪ hover_name=df_filtered.index)
      score_plot.update_layout(width= 900, height= 900,
      title="Score plot without: QCs - Negative control
      ↪ groups", font=(dict(size=16)),
      overwrite=True, legend_title=None,
      ↪ legend=dict(yanchor="top", y=-.1, xanchor="left", x=0.01, orientation="h",
      ↪ font_size=20))
      score_plot.update_traces(marker=dict(size=13))
      score_plot.add_shape(type="rect", x0=-21, x1=-14, y0=-25, y1=-20)
      score_plot.show()
```



When the negative control group, which was driving the PC1, was removed from the model the complexity of the dataset was reflected on the score plot. It was observed how the effect of digoxin on the metabolism of astrocyte cells was scattered over the two first PCs, without a clear correlation with any PC. However, the trending of the concentration was still clear, on the top left was located the no-exposure group; whereas the group exposed to the highest dose was located on the bottom right. TNF was encountered in the middle of the trending, showing that it had an impact on the metabolism similar to low-mid concentrations of digoxin; making it a suitable positive control for the experiment. Furthermore, it was observed that a sample from the no-exposure group - highlighted within a rectangle - behave as an outlier, since it was located far from the rest of the samples of its own group.

Loading plot visualization

```
[29]: loading_plot = output.loading_contourplot(pcx=0, pcy=1, histogram=True)
loading_plot.update_layout(width= 900, height= 900,overwrite=True,
                           title="Loading density plot without: QCs - Negative_
↳control groups", font=(dict(size=12)))
loading_plot
```



The loading plot kept showing an unusual behavior of the variables, and the decision to remove the positive control group as well as the outlier was taken. Hence, the PCA model would only include samples of good quality which belonged strictly to the digoxin effect.

3.5 dataset excluding: *QCs/dQCs, negative/positive control groups & outlier*

```
[30]: df_filtered = df_filtered.filter(axis=0, regex="(m)^(?!T)")

[31]: df_filtered=df_filtered.iloc[1:, :]

[32]: color = ["[C] 0.0","[C] 0.0","[C] 0.0", "[C] 0.1", "[C] 0.1",
              "[C] 0.1","[C] 0.1", "[C] 1.0", "[C] 1.0","[C] 1.
              ↪0","[C] 1.0","[C] 10.0", "[C] 10.0",
              "[C] 10.0","[C] 10.0"]

color_seq=["darkorange", "limegreen", "orchid","red"]

symbol = ["no-exposure","no-exposure",
```

```

        "no-exposure",
        ↪ "low-exposure", "low-exposure", "low-exposure", "low-exposure",
        "mid-exposure", "mid-exposure", "mid-exposure", "mid-exposure",
        ↪ "high-exposure", "high-exposure", "high-exposure",
        "high-exposure"]

symbol_seq=["diamond", "triangle-up", "triangle-down", "x", "diamond-wide"]

```

```

[33]: output_2 = unsupervised()
      pca_2 = output_2.PCA_ready(df_filtered)

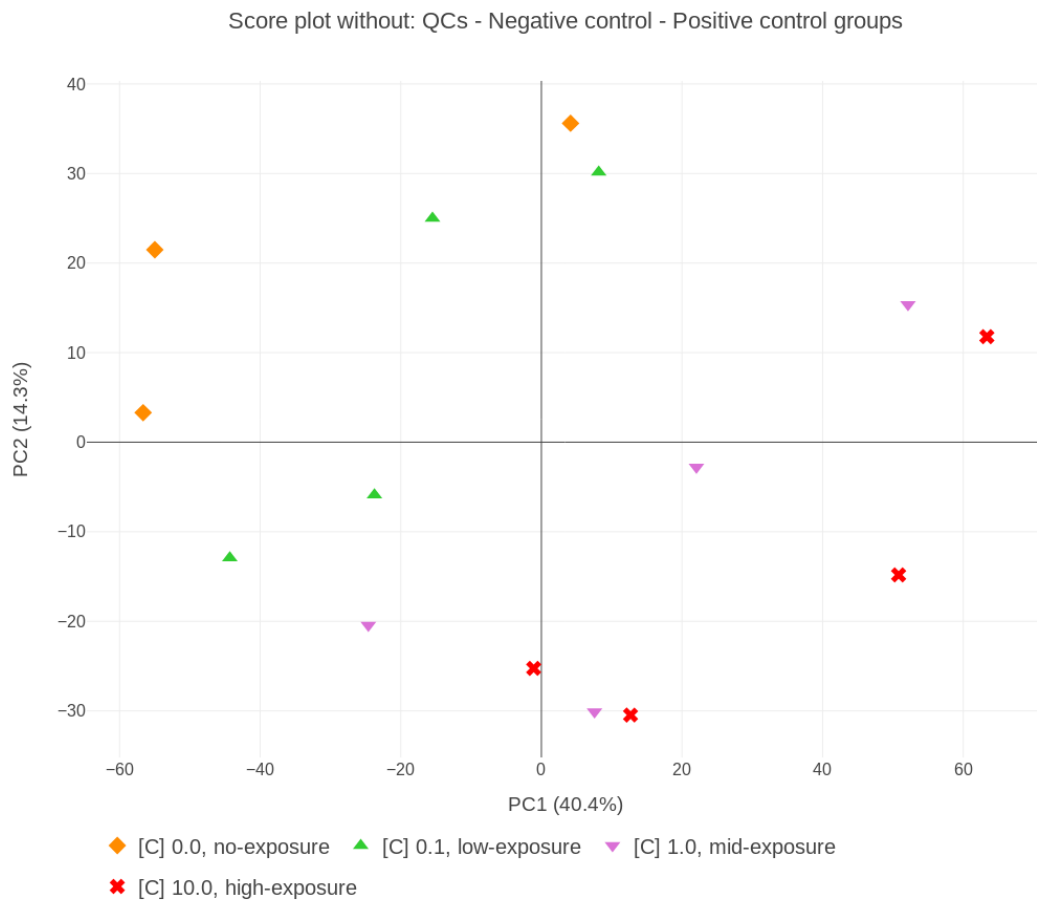
```

Score plot visualization

```

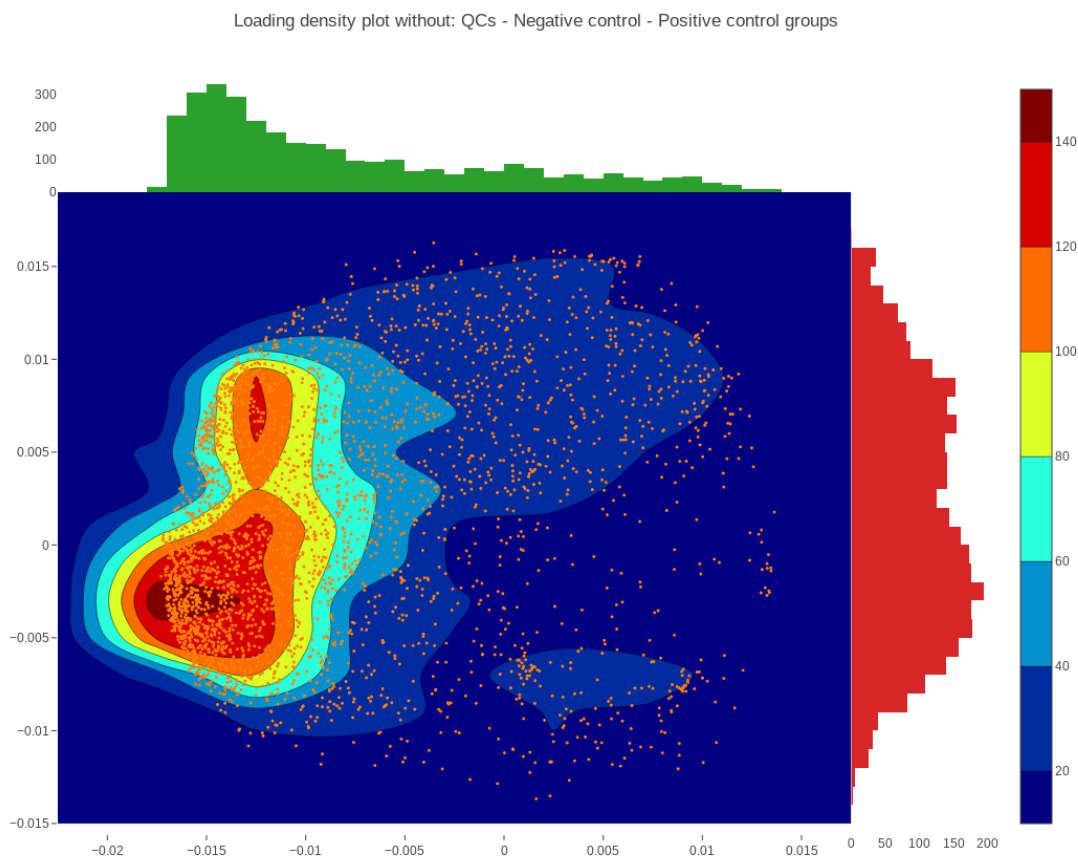
[34]: scores_plot = output_2.score_viz(pcx= 0, pcy= 1, color=color, symbol=symbol,
    ↪ label= None, label_position = None, symbol_sequence=symbol_seq,
    color_discrete_sequence=color_seq,
    ↪ hover_name=df_filtered.index)
scores_plot.update_layout(width= 900, height= 900, legend=dict(yanchor="top",
    ↪ y=-.1, xanchor="left", x=0.01, orientation="h", font_size=20),
    title="Score plot without: QCs - Negative control -
    ↪ Positive control groups", font=(dict(size=16)),
    overwrite=True, legend_title=None)
scores_plot.update_traces(marker=dict(size=13))
scores_plot.show()

```



Loading plot visualization

```
[35]: loading_plot= output_2.loading_contourplot(pcx=0, pcy=1, histogram=True)
loading_plot.update_layout(title="Loading density plot without: QCs - Negative_
control - Positive control groups", font=dict(size=12),
                           overwrite=True, width= 900, height= 900)
loading_plot
```

The last PCA model of the EDA revealed that positive control group as well as the outlier sample were masking an unusual behavior of the loadings.

On the loading plot, it was observed two hotspots of loadings. The top hotspot was considered as related to digoxin exposure, whereas the bottom hotspot seemed to have an orthogonal effect to digoxin exposure. It was hypothesized that such effect could be induced by a wide variety of possibilities, among them, it could be an uneven cell growth or mistakes on the sample preparation, etc.

Thus, it was concluded that the data was not normalized accordingly, to correct for such undesired variability. The method applied to normalize the dataset was probabilistic quotient normalization

4 Probabilistic Quotient Normalization (PQN)

Normalization is a preprocessing method which accounts for different dilutions of samples by scaling the spectra to the same virtual overall concentration. Probabilistic Quotient Normalization (PQN) was introduced as a robust normalization method more suitable for metabolomics studies compared to previous methods. PQN was based on the fact that the vast majority of analytes would not vary among samples. Then PQN computes the most probable dilution factor by looking at the

distribution of the quotients of a test sample by those of a reference profile.

4.1 Implementation

```
[36]: def PQN(path: str, sep=","):
    """
    PQN is a wrapper function for median-based Probabilistic Quotient
    ↪Normalization.

    As input, file named as "xxx1_xxx2.csv" which contains tranposed data,
    ↪where the first column is the sample names.
    Quality controls and diluted quality controls have to be named as "*QC*"
    ↪and "dQC*" respectively.

    The output is the normalised data named as "PQN_xxx1".

    Parameters
    -----
    path: path of the file to be normalised.
    sep: character separator

    References
    -----
    Frank Dieterle, Alfred Ross, Götz Schlotterbeck, Hans Senn. Probabilistic
    ↪Quotient Normalization as Robust Method to Account
    for Dilution of Complex Biological Mixtures. Application in 1H NMR
    ↪Metabonomics. Anal. Chem. 2006, 78, 4281-4290.

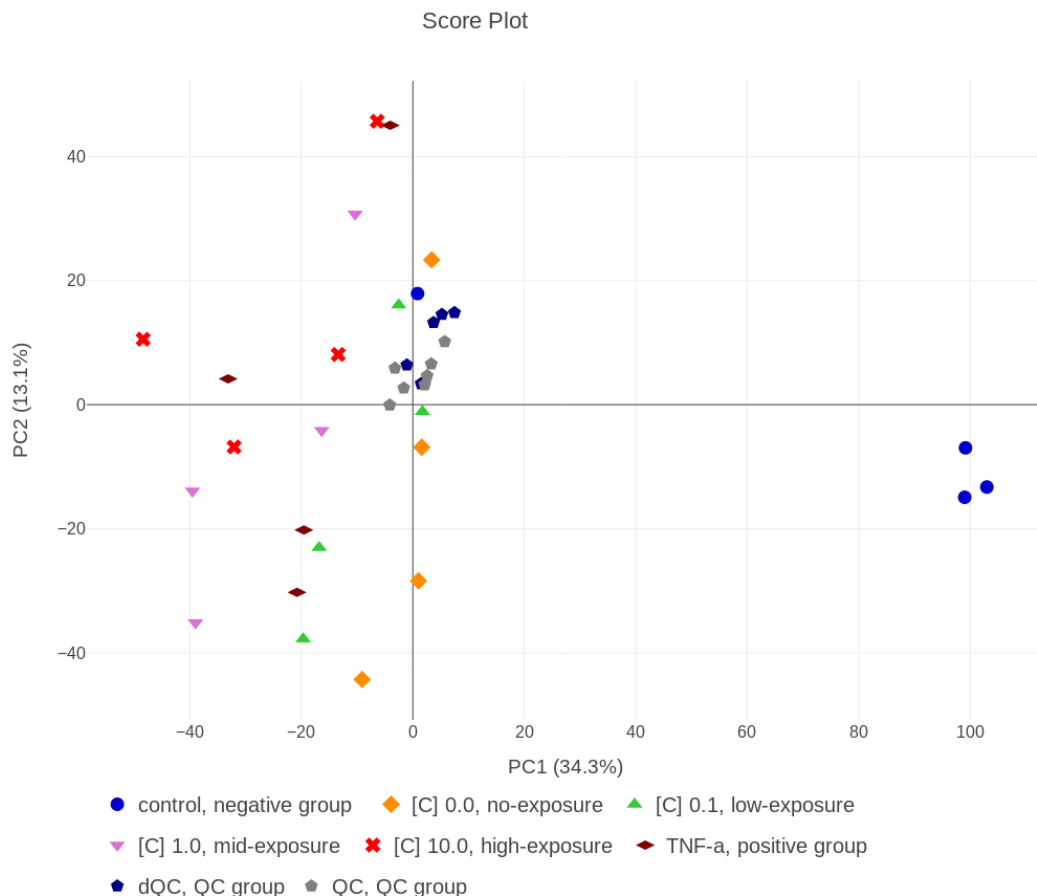
    """
    #importing the data frame from path, passing the first column as index to
    ↪use filter over it
    df = pd.read_csv(path, sep=sep, index_col=0)
    #filtering QCs and dQCs to concatenate with the normalized df later
    qcs = df.filter(axis=0, regex="^QC")
    #filtering df to obtain only QCs and computing the reference vector which
    ↪is the mean across all QCs samples
    ref_vector = np.asarray(df.filter(axis=0, regex="^QC*").median(axis=0))
    #filtering the data to exclude dQCs and QCs and using apply to compute the
    ↪quotients, variables divided by ref_vector
    quotients = df.filter(axis=0, regex="(?!QC).*").apply(lambda x: x/
    ↪ref_vector, axis=1, result_type="expand")
    #computing the coeficient vector for each sample, the scalar.
    coef_vector= np.asarray(quotients.median(axis=1))
    #filtering to exclude QC samples and computing normalization of each sample
    ↪dividing by coef_vector
    pqn_df = df.filter(axis=0, regex="(?!QC).*").apply(lambda x: x/
    ↪coef_vector, axis=0, result_type="expand")
```



```

score_plot.update_layout(width= 900, height= 900, font=dict(size=16),
    ↪legend=dict(yanchor="top", y=-.1, xanchor="left", x=0.01, orientation="h",
    ↪font_size=20),
                                overwrite=True)
score_plot.update_traces(marker=dict(size=13))
score_plot.show()

```



The suitable manner to check whether PQN was performed correctly is by having a look at the dQCs. Since dQCs are diluted samples of QCs, when applying PQN dQCs and QCs should cluster together.

This effect was observed when normalizing the data and a second EDA was carried out. However, for simplicity and readability, in this notebook it was only included the last model of EDA, containing solely digoxin exposure-related samples. By this means an easier discussion was achieved, focusing on how PQN affected the data and impaired data interpretation.

4.2 Before vs After PQN

Score plot before vs after PQN

```
[41]: df_before_pqn=df3.filter(axis=0, regex="(m)^(?!(^dQC|QC|^Control*|^T))")
```

```
[42]: color = ["[C] 0.0","[C] 0.0","[C] 0.0","[C] 0.0", "[C] 0.1", "[C] 0.1",
               "[C] 0.1","[C] 0.1", "[C] 1.0", "[C] 1.0","[C] 1.
               ↪0","[C] 1.0","[C] 10.0", "[C] 10.0",
               "[C] 10.0","[C] 10.0"]

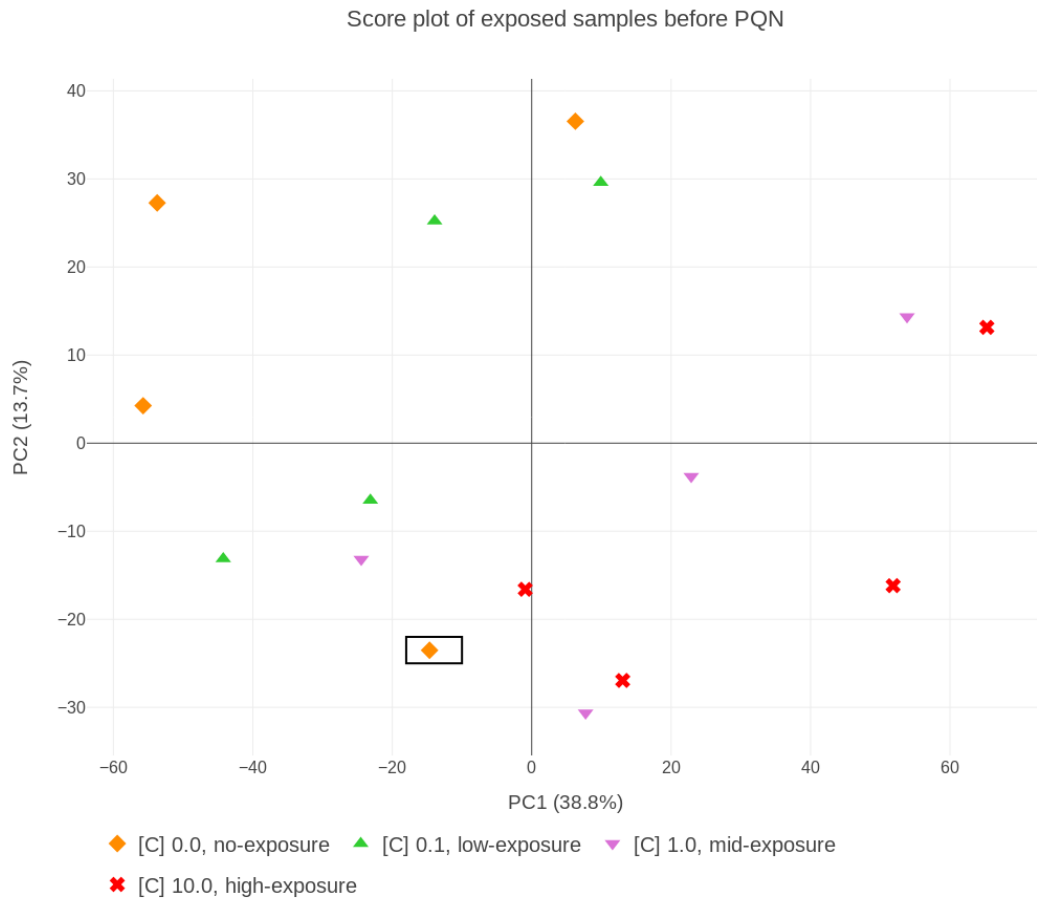
color_seq=["darkorange", "limegreen", "orchid","red"]

symbol = ["no-exposure","no-exposure","no-exposure",
          "no-exposure",↪
          ↪"low-exposure","low-exposure","low-exposure","low-exposure",
          "mid-exposure","mid-exposure","mid-exposure","mid-exposure",↪
          ↪"high-exposure","high-exposure","high-exposure",
          "high-exposure"]

symbol_seq=["diamond", "triangle-up","triangle-down", "x","diamond-wide"]
```

```
[43]: output_before = unsupervised()
pca= output_before.PCA_ready(df=df_before_pqn)
```

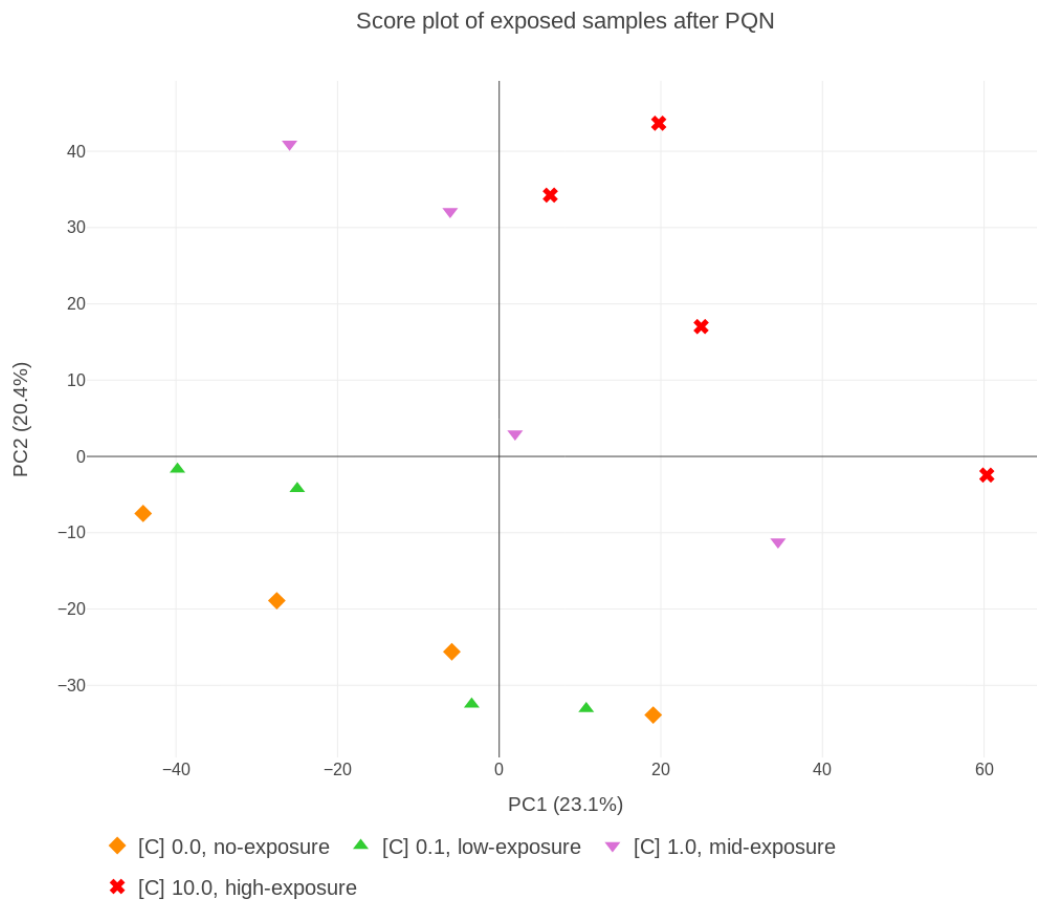
```
[44]: score_plot = output_before.score_viz(pcx= 0, pcy= 1, color=color,↪
      ↪symbol=symbol, label= None, label_position = None,↪
      ↪symbol_sequence=symbol_seq,
      color_discrete_sequence=color_seq,↪
      ↪hover_name=df_before_pqn.index)
score_plot.update_layout(width= 900, height= 900, legend=dict(yanchor="top",↪
      ↪y=-.1, xanchor="left", x=0.01, orientation="h", font_size=20),
      title="Score plot of exposed samples before PQN",↪
      ↪font=(dict(size=16)),
      overwrite=True, legend_title=None)
score_plot.update_traces(marker=dict(size=13))
score_plot.add_shape(type="rect", x0=-18, x1=-10, y0=-25, y1=-22)
score_plot.show()
```



```
[45]: df_pqn=df_pqn.filter(axis=0, regex="(m)^(?!(^dQC|QC|^Control*|^T))")
```

```
[46]: output_pqn = unsupervised()
pca= output_pqn.PCA_ready(df=df_pqn)
```

```
[47]: score_plot= output_pqn.score_viz(pcx=0, pcy=1, color=color, label=None,
    ↪label_position="top center", symbol=symbol,
    symbol_sequence=symbol_seq,
    ↪color_discrete_sequence=color_seq, hover_name=df_pqn.index)
score_plot.update_layout(legend_title=None, width= 900, height= 900,
    title="Score plot of exposed samples after PQN",
    ↪font=(dict(size=16)),
    overwrite=True, legend=dict(yanchor="top", y=-.1,
    ↪xanchor="left", x=0.01, orientation="h", font_size=20))
score_plot.update_traces(marker=dict(size=13))
score_plot.show()
```

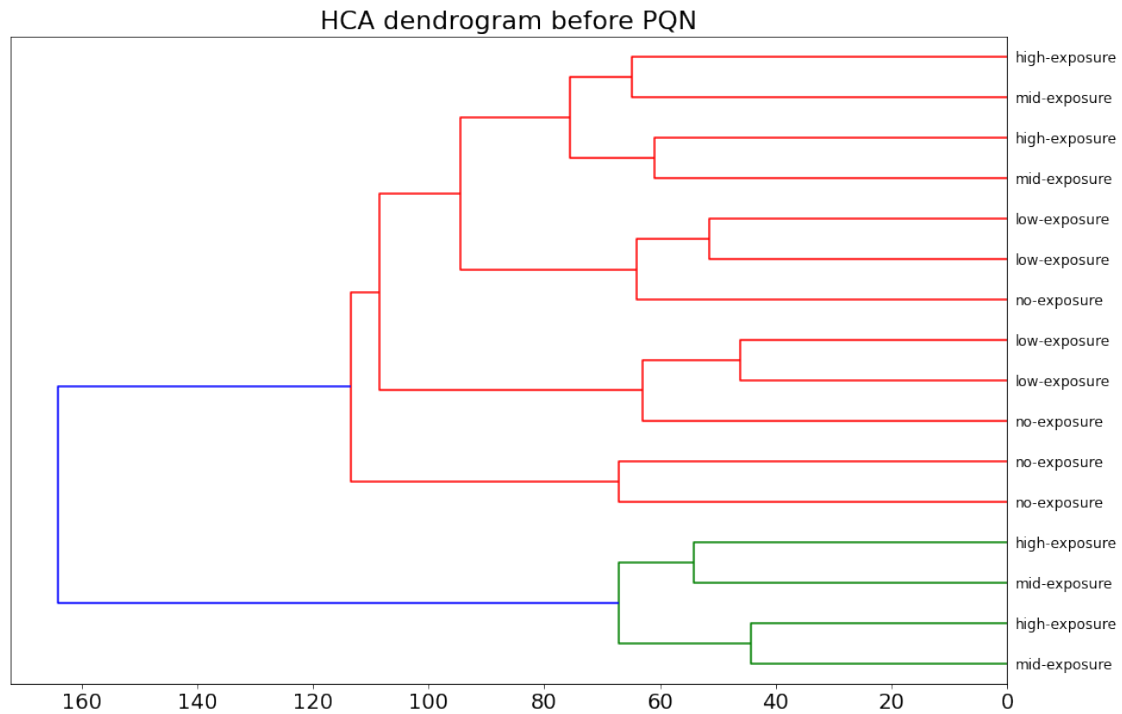


On one hand, on the score plot after PQN, it was detected no outliers; on the contrary, before applying PQN a sample belonging to no-exposure group was considered as an outlier, highlighted within a rectangle..

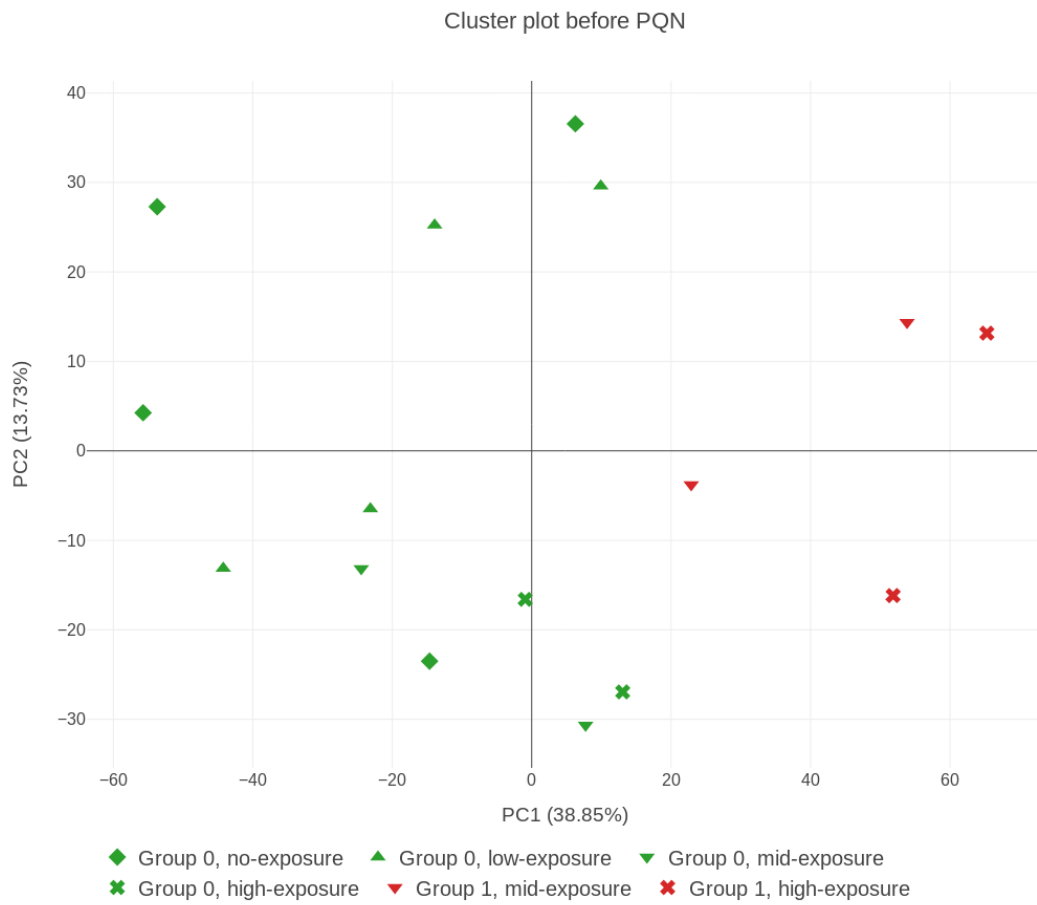
HCA before vs after PQN

```
[48]: plt.rcParams.update({'font.size': 18})
```

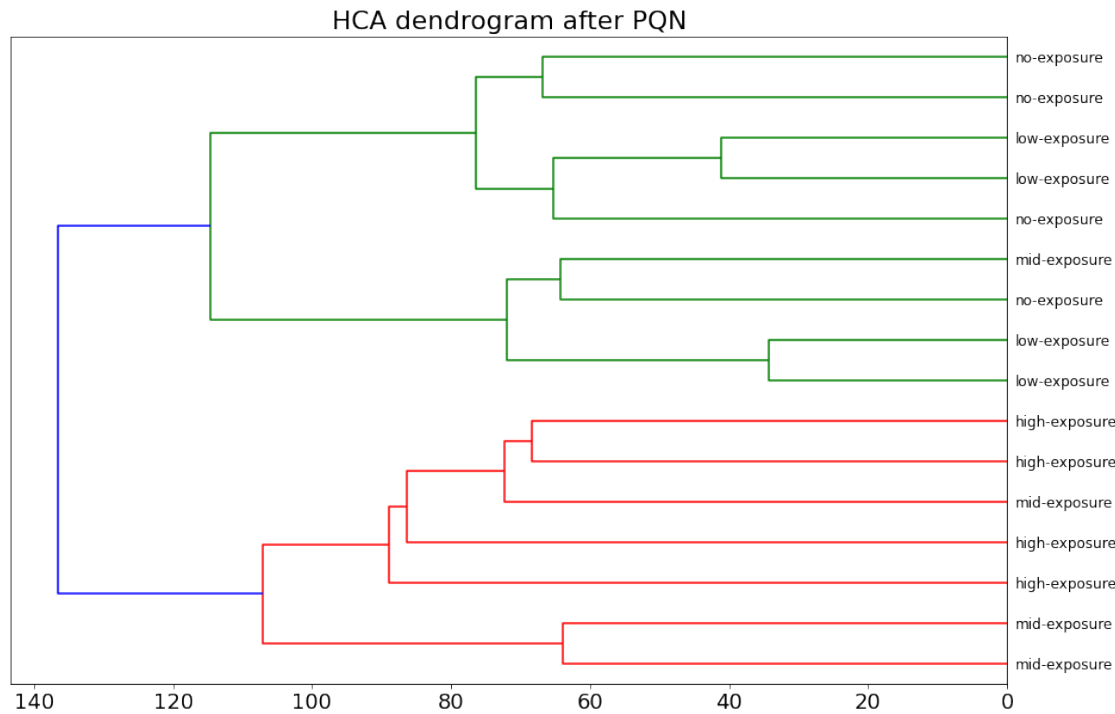
```
[49]: plt.figure(figsize=(15,10))
color_discrete_sequence=["red", "green"]
dendrogram_pca = sch.dendrogram(sch.linkage(output_before.scores,
method='ward'), labels=symbol, orientation="left")
plt.title("HCA dendrogram before PQN")
plt.show()
```



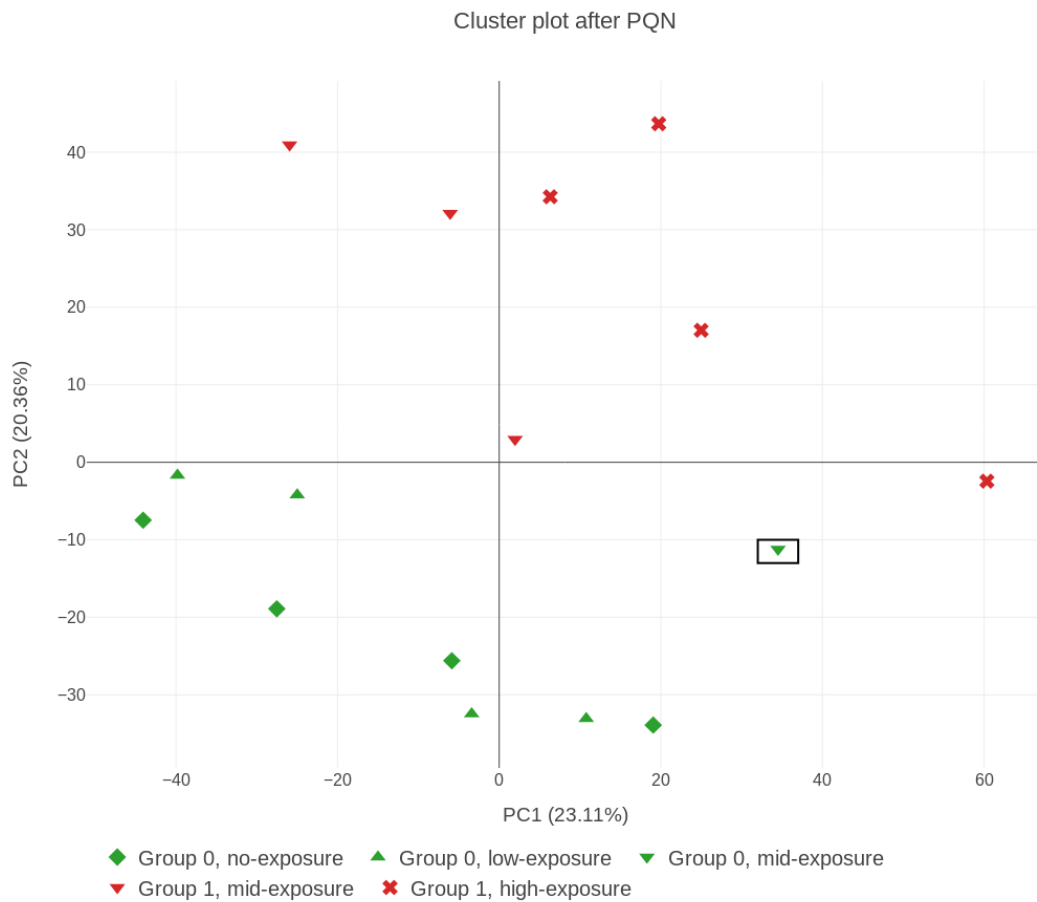
```
[50]: cluster_plot = output_before.cluster_viz(pcx=0, pcy=1, label=None,
        symbol=symbol, symbol_sequence=symbol_seq)
cluster_plot.update_layout(width= 900, height= 900, title="Cluster plot before
        PQN",legend_title=None,
        legend=dict(yanchor="top", y=-.1, xanchor="left",
        x=0.01, orientation="h",font_size=20),
        font=dict(size=16))
cluster_plot.update_traces(marker=dict(size=13))
cluster_plot.show()
```

```
[51]: plt.figure(figsize=(15,10))
sch.set_link_color_palette(["red", "green"])
dendrogram_pca = sch.dendrogram(sch.linkage(output_pqn.scores, method='ward'),
    ↳ labels=symbol, orientation="left", color_threshold=117)
plt.title("HCA dendrogram after PQN")
plt.show()
```



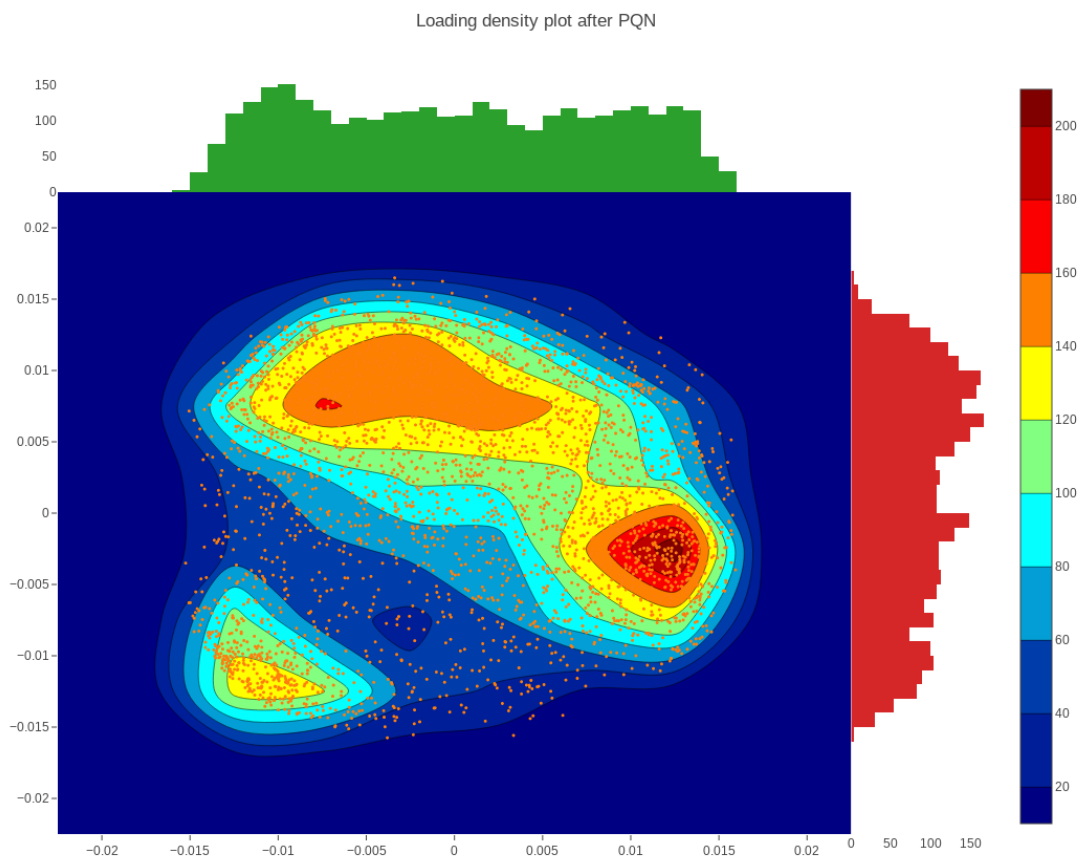
```
[52]: cluster_plot = output_pqn.cluster_viz(pcx=0, pcy=1, label=None, symbol=symbol,
      ↪symbol_sequence=symbol_seq)
cluster_plot.update_layout(width= 900, height= 900,title="Cluster plot after_
      ↪PQN", legend_title=None,
                        legend=dict(yanchor="top", y=-.1, xanchor="left",
      ↪x=0.01, orientation="h", font_size=20), font=dict(size=16))
cluster_plot.add_shape(type="rect", x0=32, x1=37, y0=-13, y1=-10)
cluster_plot.update_traces(marker=dict(size=13))
cluster_plot.show()
```



On the other hand, on the cluster plots, it was observed that before applying PQN, the HCA clustered together samples from mid/high-exposure groups with samples from the lower doses of digoxin. After applying PQN, however, HCA clustered together no/low-exposure groups separately from mid/high-exposure groups; except for a sample belonging to mid-exposure groups.

Loading plot after PQN

```
[53]: loading_plot = output_pqn.loading_contourplot(pcx=0, pcy=1, histogram=True)
loading_plot.update_layout(width= 900, height= 900, title="Loading density plot_
↳after PQN", overwrite=True, font=dict(size=12))
loading_plot
```



Through the loading plot, it was observed that there was a gradient, as now the highest concentration of loadings was related to groups with the highest concentration of digoxin. It was hypothesized that digoxin induced up-modulation of certain metabolites.

5 Supervised Analysis

The goal of supervised modelling was to obtain a robust model capable of retaining the highest amount of variability as well as presenting high prediction power. The importance of obtaining a robust model lies in the fact that it was the base to extract a subset of variables which would conform a metabolic fingerprint directly related to digoxin exposure. Since the main interest was to perform variable selection, PLS as well as RF models were chosen because one can use Variable Important in Projection (VIP) and Variable Importance, respectively, to carry out that task.

Within supervised modelling process, both frameworks can be applied, either classification or regression. A gold-standard model in metabolomics is PLS-DA, a classification model; however in this section it is discussed the limitations observed when used the classification framework on our metabolomic data, which presented highly unbalance structure with $p \gg n$.

5.1 Classification Framework

Even though, neither PLS-DA nor RF-classifier performed good enough to be retained as final models, here it was addressed one of the hardest task when using a RF-classifier model, its interpretability.

For the classification framework, the response y variable was defined as discrete, referring the concentration groups, now classes, which were no-exposure, low-exposure, mid-exposure and high-exposure. The response variable was encoded to obtain a unitary dummy matrix.

```
[54]: np.random.seed(42)
```

```
[55]: df = pd.read_csv("digoxin_pqn.csv", index_col=0)
df = df.filter(axis=0, regex="(m)^(?!(^dQC|QC|^Control*|^T).)")
```

```
[56]: exposure = ["exposed_0.0", "exposed_0.0", "exposed_0.0", "exposed_0.0",
                  "exposed_0.1", "exposed_0.1", "exposed_0.1", "exposed_0.1",
                  ↪ "exposed_1.0", "exposed_1.0", "exposed_1.0", "exposed_1.0",
                  "exposed_10.0", "exposed_10.0", "exposed_10.0", "exposed_10.0"]
```

```
[57]: color = ["[C] 0.0", "[C] 0.0", "[C] 0.0", "[C] 0.0", "[C] 0.1", "[C] 0.1",
               "[C] 0.1", "[C] 0.1", "[C] 1.0", "[C] 1.0", "[C] 1.0",
               ↪ "[C] 1.0", "[C] 10.0", "[C] 10.0", "[C] 10.0", "[C] 10.0", "[C] 10.0"]

color_seq=["darkorange", "limegreen", "orchid", "red"]
```

```
symbol = ["no-exposure", "no-exposure", "no-exposure", "no-exposure",
          ↪ "low-exposure", "low-exposure", "low-exposure", "low-exposure",
          "mid-exposure", "mid-exposure", "mid-exposure", "mid-exposure",
          ↪ "high-exposure", "high-exposure", "high-exposure", "high-exposure"]

symbol_seq=["diamond", "triangle-up", "triangle-down", "x", "diamond-wide"]
```

```
[58]: class_group=["no-exposure", "low-exposure", "mid-exposure", "high-exposure"]
```

```
[59]: from sklearn.preprocessing import LabelBinarizer
from chemometrics.supervised import Supervised
```

```
[60]: encoder = LabelBinarizer()
y_classes = encoder.fit_transform(exposure)
y_classes
```

```
[60]: array([[1, 0, 0, 0],
            [1, 0, 0, 0],
            [1, 0, 0, 0],
```

```
[1, 0, 0, 0],
[0, 1, 0, 0],
[0, 1, 0, 0],
[0, 1, 0, 0],
[0, 1, 0, 0],
[0, 0, 1, 0],
[0, 0, 1, 0],
[0, 0, 1, 0],
[0, 0, 1, 0],
[0, 0, 0, 1],
[0, 0, 0, 1],
[0, 0, 0, 1],
[0, 0, 0, 1]]])
```

5.1.1 PLS-DA

```
[61]: output = Supervised()
```

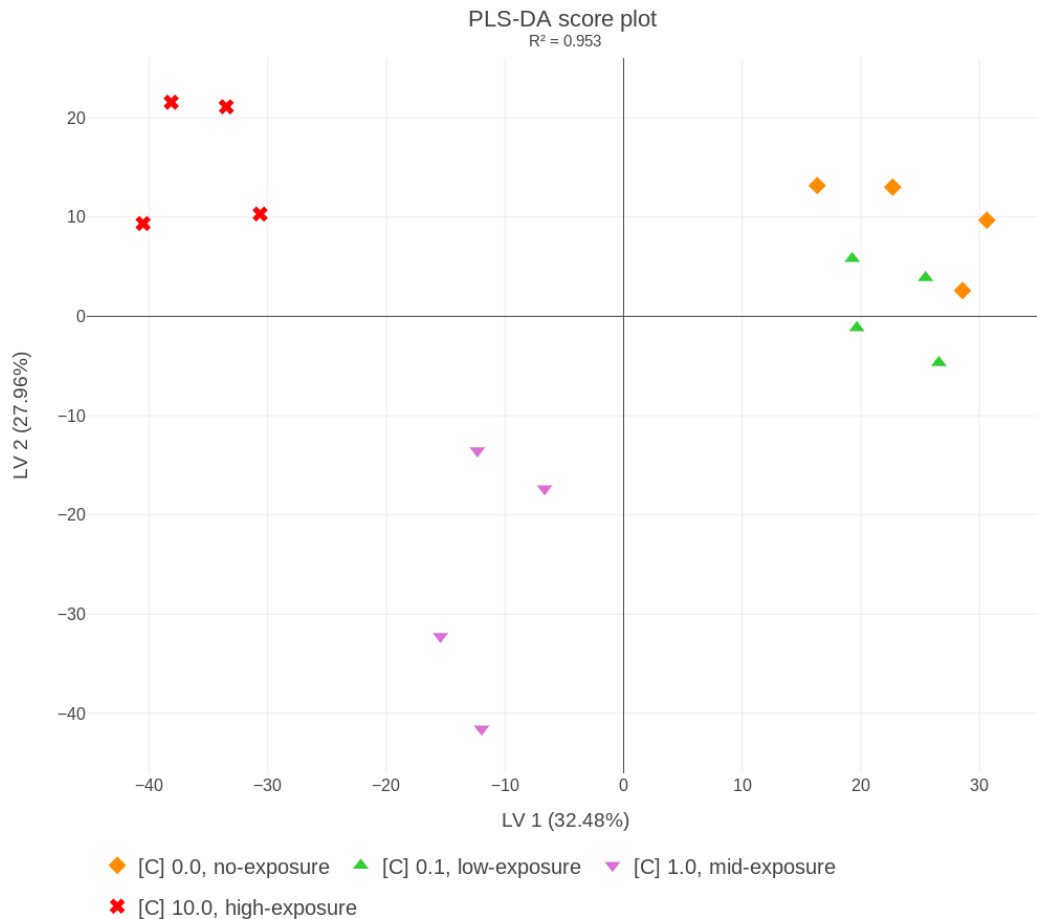
```
[62]: pls_da = output.PLS_model(df, y_classes, n_components=4)
```

```
[63]: print(f'Coefficient of determinarion R2 for Y = {output.pls_coef_determ:.3f}')
```

Coefficient of determinarion R² for Y = 0.953

Score plot visualization

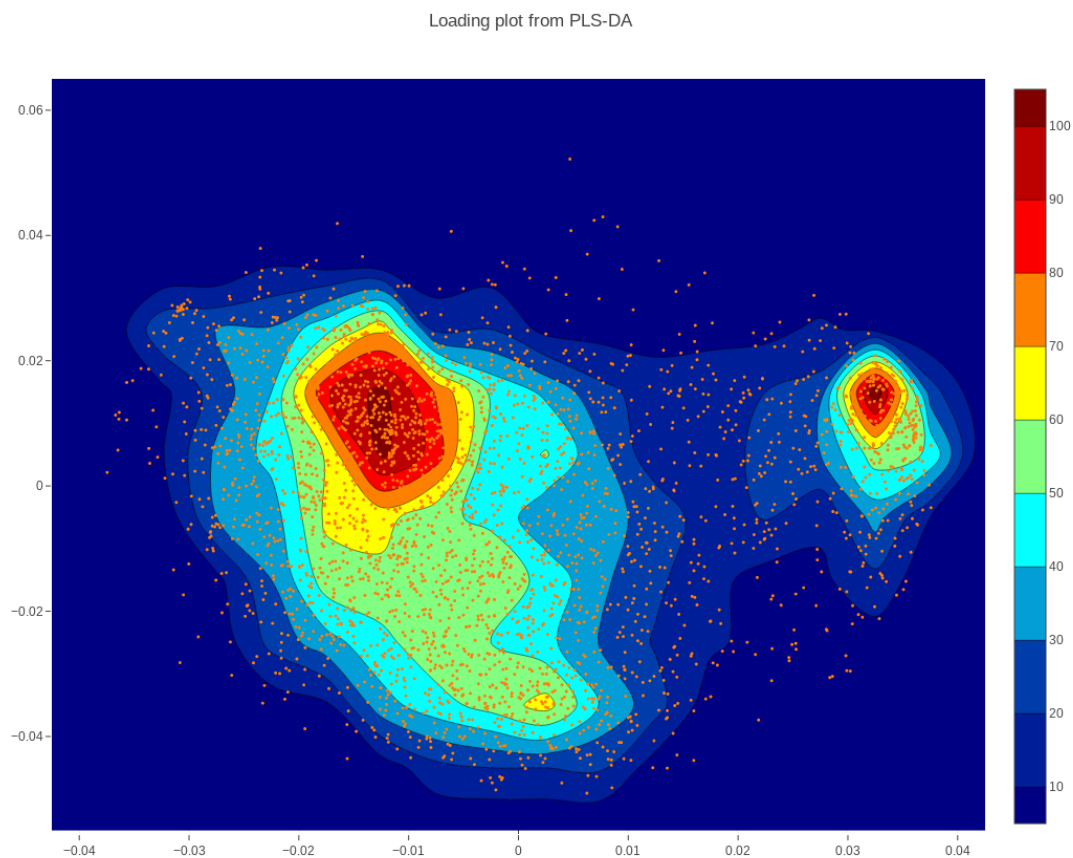
```
[64]: pls_da_scoreplot = output.Score_viz(lv_1=0, lv_2=1,name_list=None,
    ↪symbol=symbol, symbol_sequence=symbol_seq, color=color,
    ↪color_discrete_sequence=color_seq)
pls_da_scoreplot.update_layout(title="PLS-DA score plot <br><sup>R2 = 0.953</sup>"/
    ↪sup>",
                                legend=dict(yanchor="top", y=-.1,
    ↪xanchor="left", x=0.01, orientation="h",font_size=20),
                                width= 900, height= 900, font=(dict(size=16)),
    ↪legend_title=None, overwrite=True)
pls_da_scoreplot.update_traces(marker=dict(size=13))
pls_da_scoreplot.show()
```



The PLS-DA model reported a $R^2 = 0.953$, and it managed to discriminate between the groups of high-exposure and no/low-exposure in the first Latent Variable (LV), localizing the mid-exposure group somewhere in the middle between the extreme groups

Loading plot visualization

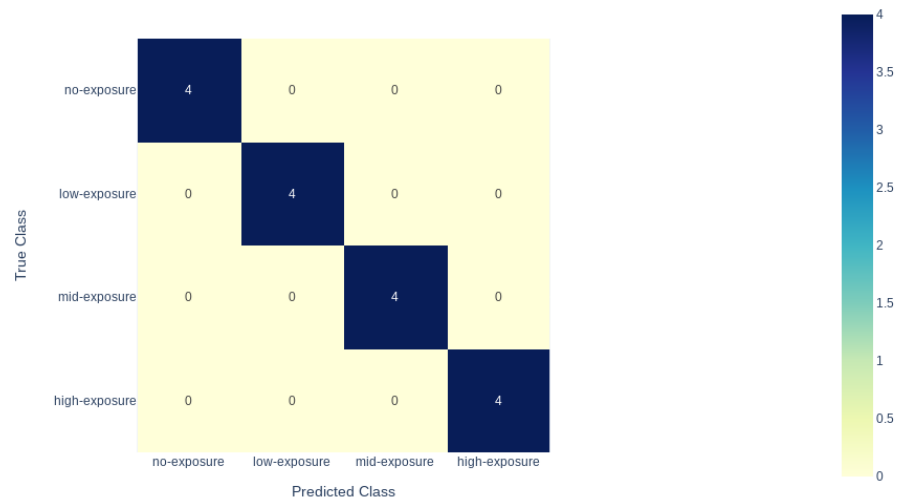
```
[65]: pls_da_loadingplot = go.Figure()
pls_da_loadingplot.add_trace(go.Histogram2dContour(x=pls_da.x_loadings[:,0],
    ↪ y=pls_da.x_loadings[:,1], colorscale="Jet", xaxis="x", yaxis="y"))
pls_da_loadingplot.add_trace(go.Scatter(x=pls_da.x_loadings[:,0], y=pls_da.
    ↪ x_loadings[:,1], mode="markers", marker=dict(size=3), xaxis="x", yaxis="y"))
pls_da_loadingplot.update_layout(template="presentation", width= 900, height=
    ↪ 900, hovermode=False, font=dict(size=12), title= "Loading plot from PLS-DA")
pls_da_loadingplot.show()
```



Classification performance metrics

```
[66]: from chemometrics.utils import ConfusionMat_viz, ROC_AUC_viz
```

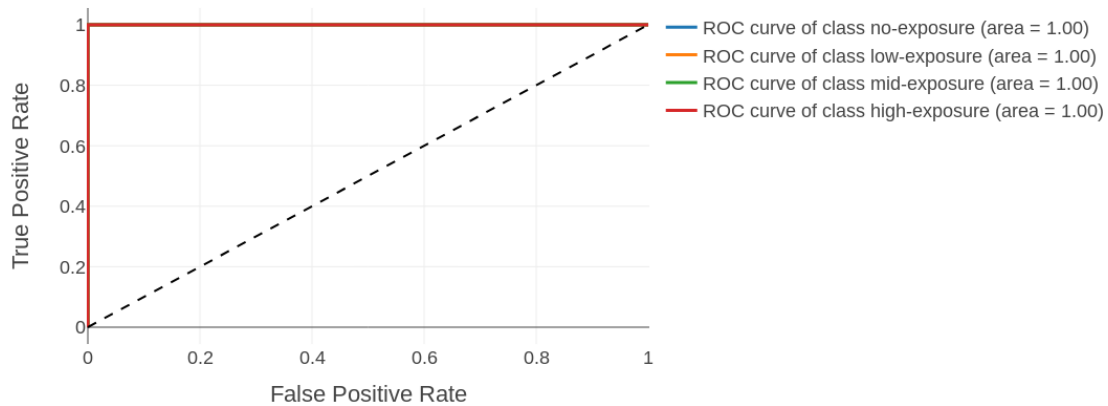
```
[67]: ConfusionMat_viz(y=y_classes, y_hat=output.y_pred, classes_names=class_group)
```

```
[68]: output.Classification_report(y_variable=y_classes, y_hat=output.y_pred,
    ↪target_names=class_group)
```

	precision	recall	f1-score	support
no-exposure	1.00	1.00	1.00	4
low-exposure	1.00	1.00	1.00	4
mid-exposure	1.00	1.00	1.00	4
high-exposure	1.00	1.00	1.00	4
accuracy			1.00	16
macro avg	1.00	1.00	1.00	16
weighted avg	1.00	1.00	1.00	16

```
[69]: ROC_AUC_viz(n_classes=class_group, y_classes=y_classes, y_hat=output.y_pred)
```



It was observed a perfect classification from PLS-DA, where no single sample was misclassified.

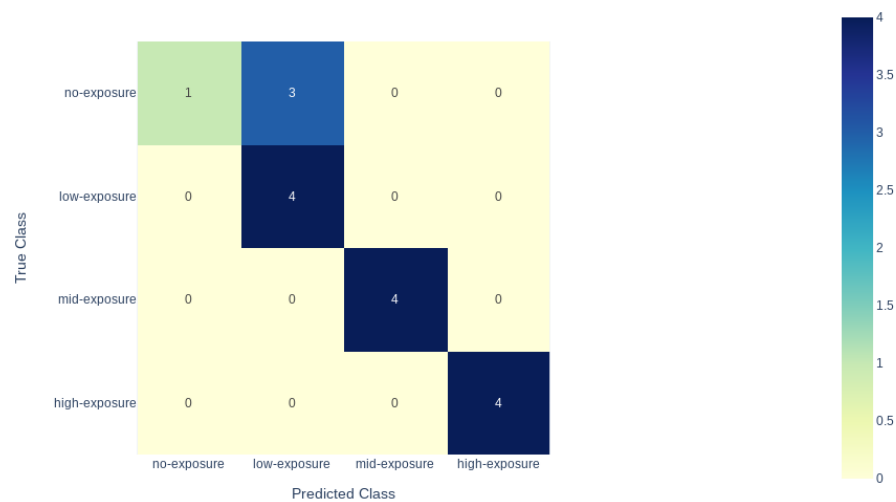
Cross-Validation

```
[70]: from sklearn.model_selection import cross_val_predict, LeaveOneOut
```

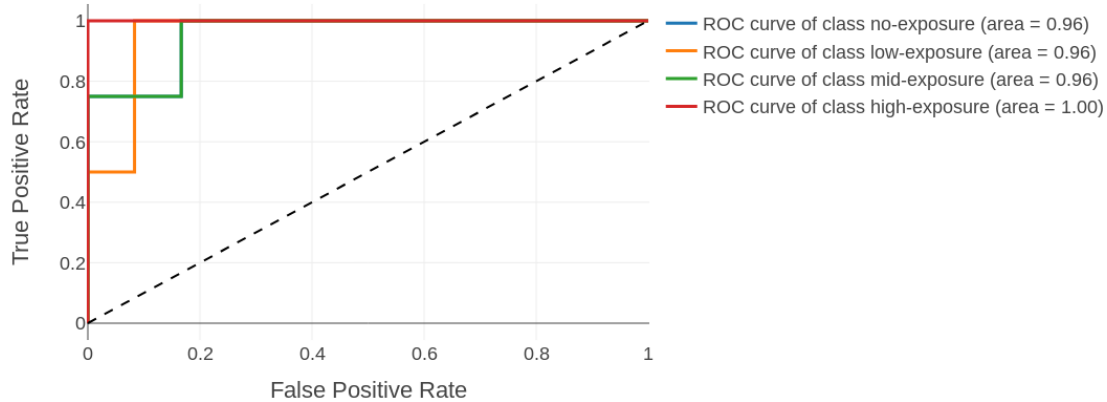
```
[71]: cv = LeaveOneOut()
```

```
[72]: pls_da_cv_pred = cross_val_predict(estimator=pls_da, X=df, y=y_classes, cv=cv)
```

```
[73]: ConfusionMat_viz(y= y_classes, y_hat=pls_da_cv_pred, classes_names=class_group)
```



```
[74]: ROC_AUC_viz(n_classes=class_group, y_classes=y_classes, y_hat=pls_da_cv_pred)
```



```
[75]: output.Classification_report(y_variable=y_classes, y_hat=pls_da_cv_pred,
    ↪target_names=class_group)
```

	precision	recall	f1-score	support
no-exposure	1.00	0.25	0.40	4
low-exposure	0.57	1.00	0.73	4
mid-exposure	1.00	1.00	1.00	4
high-exposure	1.00	1.00	1.00	4
accuracy			0.81	16
macro avg	0.89	0.81	0.78	16
weighted avg	0.89	0.81	0.78	16

Quoting scikit-learn:

“Precision is the ability of the classifier no to label as positive a sample that is negative and recall is the ability of the classifier to find all positive samples.

The F-1 measure can be interpreted as a weighted harmonic mean of the precision and recall. A F-1 measure reaches its best value at 1 and it worst score at 0.

In a binary classification problem, the terms “positive” and “negative” refer to the classifier’s prediction, and terms *true* and *false* refer to whether that prediction corresponds to the observation. From a confusion matrix precision, recall and F-1 are derived as follows:

- $precision = \frac{TP}{TP+FP}$

- $recall = \frac{TP}{TP+FN}$
- $F-1 = 2 * \frac{precision*recall}{precision+recall}$

”

In this context, it was observed that PLS-DA failed to correctly predict samples belonging to **no-exposure** class, returning a high number of FN with a $recall = 0.25$.

The FN samples were misclassified as **low-exposure**; for this reason $precision = 0.57$ and the class **low-exposure** has a high number of FP.

Model evaluation and validation The aim of CV was to compare the goodness of fit R^2 and goodness of prediction Q^2 . R^2 , also known as coefficient of determination, represents the fraction of variability which the model is capable of explaining, whereas the Q^2 represents the prediction capability of the model. Goodness of fit and goodness of prediction metrics involves Total Sum of Squares (TSS), Residual Sums of Squares (RSS) as well as Predicted Error Sum of Squares (PRESS):

- $TSS = \sum (y_i - \bar{y})^2$
- $RSS = \sum (y_i - \hat{y})^2$
- $PRESS = \sum (y_i - \hat{y}_i)^2$

Where the difference between R^2 and Q^2 is that PRESS is computed during CV, when observations are kept out:

- $R^2 = 1 - \frac{RSS}{TSS}$
- $Q^2 = 1 - \frac{PRESS}{TSS}$

```
[76]: output.PLS_goodness_metrics(X=df, y_variable=y_classes, n_models=4)
```

```
[77]: goodness_met_viz = output.PLS_goodness_metrics_viz()
goodness_met_viz.update_layout(title="PLS-DA goodness metrics",
                               width= 400, height= 700, font=(dict(size=12)),
                               ↪overwrite=True)
goodness_met_viz.update_xaxes(type="category")
goodness_met_viz.show()
```



When the goodness metrics were evaluated, the goodness of prediction behaved extremely unusual. A normal behavior of Q² is that it starts to notably increase when the model includes 1 and 2 LV and afterwards it starts to reach a plateau or slightly decrease. For this reason, PLS-DA failed model validation and was rejected.

5.1.2 RF Classifier

```
[78]: output_rf_c = Supervised()
```

```
[79]: rf_c = output_rf_c.RandomForest_Classifier(X=df, y_variable=y_classes)
```

```
[80]: print(f'OOB score = {rf_c.oob_score_}')
```

OOB score = 0.0625

Classification performance metrics

```
[81]: y_pred_rf_c = rf_c.predict(df)
```

```
[82]: output_rf_c.Classification_report(y_variable=y_classes, y_hat=y_pred_rf_c,
    ↪ target_names=class_group)
```

	precision	recall	f1-score	support
no-exposure	1.00	1.00	1.00	4

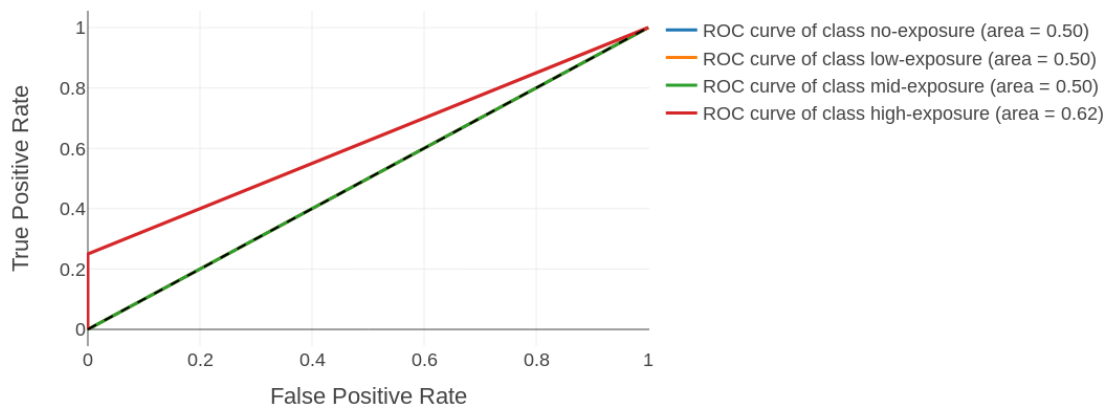
low-exposure	1.00	1.00	1.00	4
mid-exposure	1.00	1.00	1.00	4
high-exposure	1.00	1.00	1.00	4
accuracy			1.00	16
macro avg	1.00	1.00	1.00	16
weighted avg	1.00	1.00	1.00	16

```
[83]: print(f'OOB score = {rf_c.oob_score_}')
```

OOB score = 0.0625

OOB score is really low, that means that the model had poor prediction power and that must be due to the low number of observations present in the dataset. For that reason, even though RF managed to correctly classify each sample, it failed in prediction when CV was applied.

```
[84]: ROC_AUC_viz(n_classes=class_group, y_classes=y_classes,
    ↪ y_hat=cross_val_predict(rf_c, X=df, y=y_classes, cv=cv))
```



ROC curve showed that the RF-classifier would predict half of the time one class for each samples.

```
[85]: import seaborn as sns
```

```
[86]: def distanceMatrix(model, X):
    terminals = model.apply(X)
    #[n_samples, n_trees] containing the terminal node for each sample in each
    ↪ tree
    nTrees = terminals.shape[1]
```

```

#getting the nodes of the samples in the first tree
a = terminals[:,0]
#assess whether values are equal or not and by multiplying by 1 it obtained
↪ a binary matrix
#initializing the first freq count
proxMat = 1*np.equal.outer(a, a)

for i in range(1, nTrees):
    a = terminals[:,i]
    proxMat += 1*np.equal.outer(a, a)

distanceMat = np.absolute(np.round(1-(proxMat / nTrees), decimals=1))

return distanceMat

```

```

[87]: rf_c_distMat = distanceMatrix(rf_c, df)
      rf_c_distMat=pd.DataFrame(rf_c_distMat)

```

```

[88]: rf_c_distMat.set_axis(df.index, axis=1, inplace=True)
      rf_c_distMat.set_axis(df.index, axis=0, inplace=True)

```

5.2 Regression Framework

For the regression framework, the response variable y was defined as the concentration to which the samples were exposed ranging from 0.0 to 10 μ molar.

The PLS-RA model reported an $R^2 = 0.999$ which outperformed both PLS-DA as well as RF regressor which reported an $R^2 = 0.955$ with an OOB score still showing lack of prediction power, $OOB = 0.570$.

```

[89]: y_cont = np.array([0.0, 0.0,0.0,0.0,0.1,0.1,0.1,0.1,1.0,1.0,1.0,1.0,10.0,10.
↪ 0,10.0,10.0])

```

5.2.1 PLS-RA

```

[90]: output_2 = Supervised()
      pls_ra = output_2.PLS_model(X=df, y_variable=y_cont, n_components=4)

```

```

[91]: print(f' Coefficient of determination  $R^2$  for Y = {output_2.pls_coef_determ:.
↪ 3f}')

```

Coefficient of determination R^2 for Y = 0.999

Score plot visualization

```

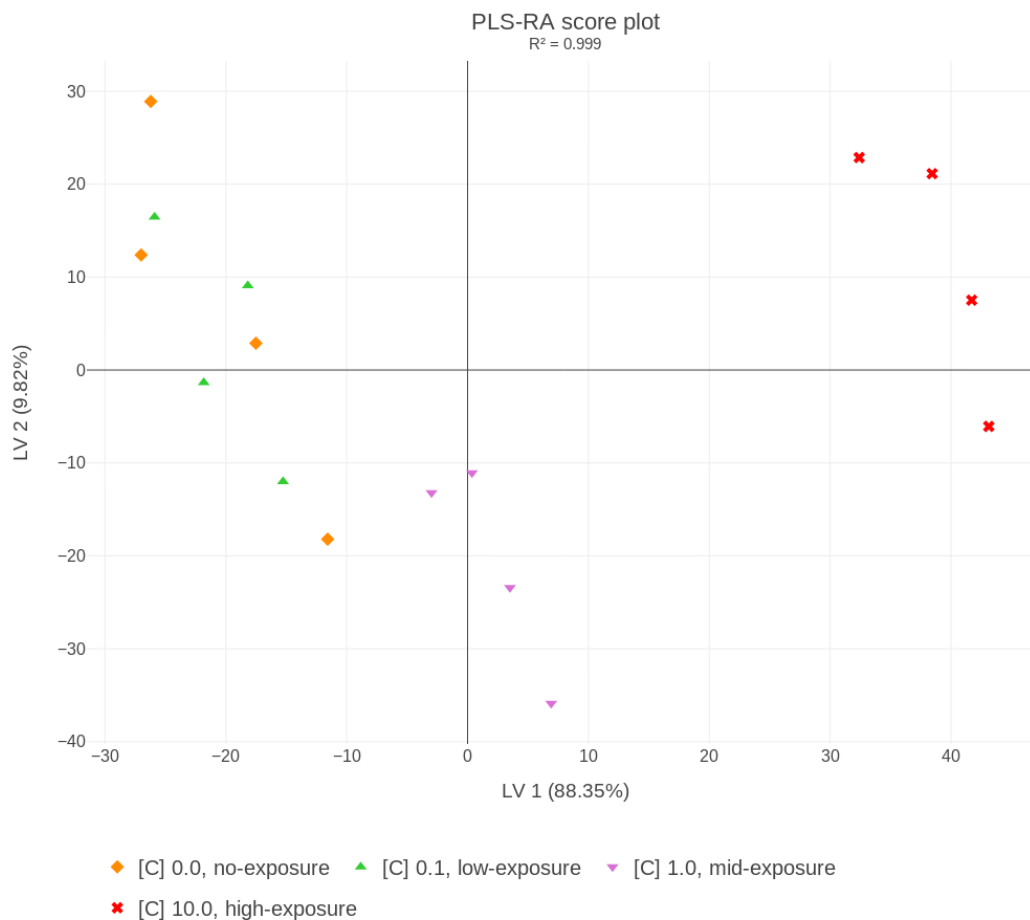
[92]: score_plot_2 = output_2.Score_viz(lv_1=0, lv_2=1, name_list=None, color=color,
↪ symbol=symbol, symbol_sequence=symbol_seq, color_discrete_sequence=color_seq)

```

```

score_plot_2.update_layout(title="PLS-RA score plot<br><sup>R2 = 0.999",
    ↪legend=dict(yanchor="top", y=-.15, xanchor="left", x=0.01, orientation="h",
    ↪font_size=20),
    width= 900, height= 900, font=(dict(size=16)),
    ↪overwrite=True)
score_plot_2.update_traces(marker=dict(size=10))
score_plot_2.show()

```



Strikingly, PLS-RA score plot, showed harmony with the classification result from RF-classifier; showing that the group of high exposure was the more dissimilar to the rest. Such situation was retained in the first LV, where on the left-hand side rested no/low-exposure groups and, extremely close to them, the mid-exposure group. On the right-hand side, however, it was located the high-exposure group.

Model evaluation and validation

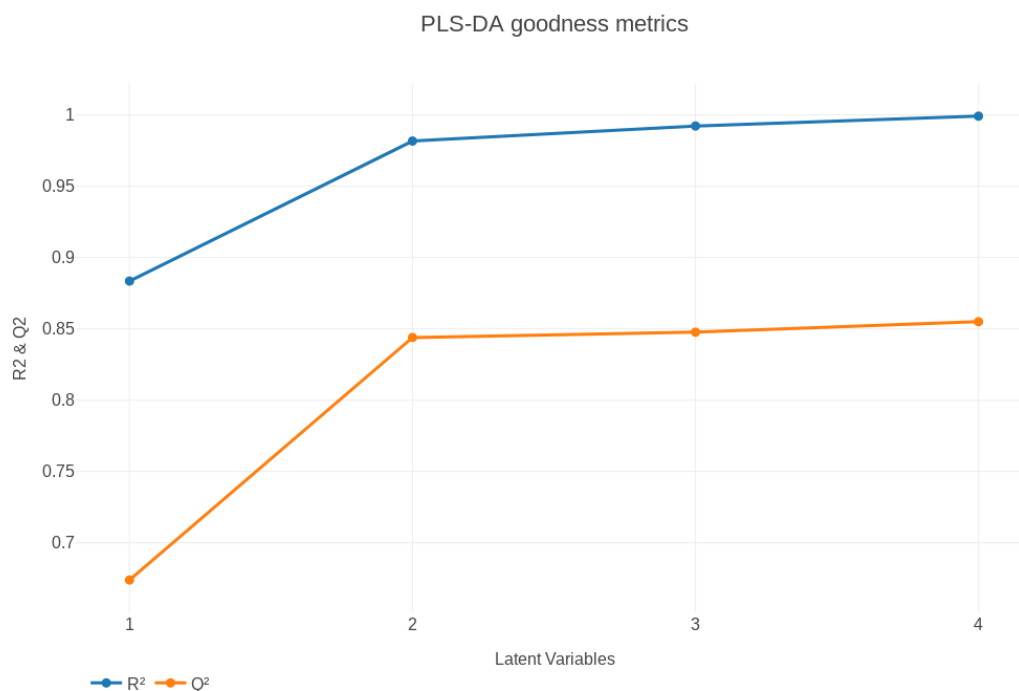
```
[93]: output_2.PLS_goodness_metrics(X=df, y_variable=y_cont, n_models=4)
```

```
[94]: output_2.r2_nested
```



```
[94]: [0.8834663026277807, 0.9816321692364528, 0.9922095874067319, 0.999084079879463]
```

```
[95]: goodness_met_viz=output_2.PLS_goodness_metrics_viz()
goodness_met_viz.update_layout(title="PLS-DA goodness metrics",
                               legend=dict(orientation="h", yanchor="top", y=-.
↪10, xanchor="left", x=0.01),
                               width= 400, height= 700, font=(dict(size=16)),↪
↪overwrite=True)
goodness_met_viz.update_xaxes(type="category")
goodness_met_viz.show()
```



PLS-RA goodness metrics were satisfactory. It was observed that a model including only two LV would perform well in terms of explained variability as well as prediction power; thus, accomplishing the parsimony criteria of the simplest model.

As shown, the PLS-RA model successfully passed model validation, for that reason it was retained and further explored to extract a metabolic fingerprint. The next step consisted on variable selection to obtain the metabolic fingerprint. Variable selection was done by computation of Variable Importance in Projection (VIP) of the PLS-RA model.

6 Metabolic Fingerprint Extraction

As aforementioned, PLS model was chosen since it allows variable selection, which means that it allows the extraction of a subset of variables of interest for further study. In the PLS framework,

there is a wide variety of methods to perform variable selection, in this project it was use Variable Importance in Projection (VIP) which belongs to the category of filtering methods.

The idea behind VIP is to accumulate the importance of each variable being reflected by loading weights from each component. In the scikit-learn python library the computation of the VIP was not implemented yet, hence VIP were computed by the following the equation from where VIP measure, v_j , is defined by the following equation:

$$v_j = \sqrt{p \sum_{a=1}^A [(q_a^2 t_a' t_a) (W_{aj} / \|W_a\|^2)] / \sum_{a=1}^A (q_a^2 t_a' t_a)}$$

Implemented as follows:

```
[96]: def VIP(x, model):
    """
    Computes PLS Variable Importance in Projection (VIP)

    Parameter
    -----
    x: data frames
    model: scikit-learn PLS model

    Return
    -----
    VIP

    References
    -----
    Tahir Mehmood, Kristian Hovde Liland, Lars Snipen, and Solve Sæbø. A review
    ↪ of variable
    ↪ selection methods in partial least squares regression. Chemometrics and
    ↪ intelligent laboratory
    ↪ systems, 118:62-69, 2012

    Author
    -----
    https://github.com/scikit-learn/scikit-learn/issues/7050
    """
    t = model.x_scores_
    w = model.x_weights_
    q = model.y_loadings_

    m, p = x.shape
    _, h = t.shape

    vips = np.zeros((p,))

    s = np.diag(t.T @ t @ q.T @ q).reshape(h, -1)
    total_s = np.sum(s)
```

```

    for i in range(p):
        weight = np.array([ (w[i,j] / np.linalg.norm(w[:,j]))**2 for j in
↪range(h) ])
        vips[i] = np.sqrt(p*(s.T @ weight)/total_s)

    return vips

```

```
[97]: vip = VIP(x=df, model=pls_ra)
```

Variable filtering was done based on three different thresholds which were tested to evaluate how it affected the performance of the PLS-RA model and how many informative variables could be extracted. For a $VIP_{threshold} = X$, variables were retained if their $VIP_{value} \geq VIP_{threshold}$. New PLS-RA models were built and its goodness metrics were studied to assess how the reduction in the number of variables affected the models.

The following function was developed to automatize VIP **subsetting**:

```
[98]: def VIP_subset(vips, threshold, X):
    """
    Performs df extraction from original df according to the threshold value
↪passed

    Parameter
    -----
    - vips: data frame or data matrix containing VIP values
    - threshold: threshold value, it can be integer or float
    - X: data frame containing annotated metabolites from where extraction is
↪carried out based on the VIP threshold and values

    Return
    -----
    Data frame containing the extracted metabolites of interest. The content is
↪the values of the original data frame

    Author:
    -----
    Christian Peralta

    """
    # vip matrix is converted to a pd.DataFrame and transpose, so later the
↪variable name can be included in the output
    v = pd.DataFrame(vips).T
    #filtering according to a selected threshold
    v = v[v>=threshold].dropna(axis=1)
    # DF-v now contains the variables of interest according to the VIP
↪threshold.

```

```

# Variable name is set to object and stored in a variable for iteration
features_subset = v.columns.astype("object")
# setting an empty dictionary which will be converted to a DF
d = {}
# loop setting as key:value argument as variable_position:column of the DF-X
for i in features_subset:
    d[i] = X[str(i)]
d = pd.DataFrame(d)
return d

```

```
[99]: var_subset = VIP_subset(vips=vip,threshold=2, X=df)
```

The following functions was developed to automatically perform VIP feature selection which then in pass to PLS for model fitting. The main goal was to automatically compute goodness metrics for nested models at desired $VIP_{thresholds}$.

```
[100]: def comparative_goodness_metrics(threshold=None):
    """
    Computes goodness metrics for specified VIP subsets, allowing the
    ↪comparison in terms of nested  $R^2$  and  $Q^2$ .
    It takes as parameter treshold for VIP_subset function

    Parameter
    -----
    - threshold: integer value set as VIP threshold

    Return
    -----
    List of dictionaries containing  $R^2$ ,  $Q^2$  values as well as beta-coefficients
    ↪for the corresponding VIP_threshold

    Author
    -----
    Christian Peralta
    """

    # first, variable subset extraction based on VIP threshold, using the
    ↪aforedeclared function 'VIP_subset'
    # subset DF will be use to fit PLS nested models
    subset = VIP_subset(vips = vip, threshold=threshold, X=df)
    # fitting and computation of goodness metrics
    results = output_2.PLS_goodness_metrics(X=subset, y_variable=y_cont,
    ↪n_models=4, vip_metrics=True, treshold=threshold)
    return results

```

```
[101]: # using map to iteratively compute metrics of the 3 desired thresholds
tr_1, tr_1_5, tr_2 = list(map(comparative_goodness_metrics, [1,1.5,2]))
```

```
/home/christian97/.local/lib/python3.8/site-  
packages/sklearn/cross_decomposition/_pls.py:507: FutureWarning:
```

The attribute ``coef_`` will be transposed in version 1.3 to be consistent with other linear models in scikit-learn. Currently, ``coef_`` has a shape of `(n_features, n_targets)` and in the future it will have a shape of `(n_targets, n_features)`.

```
/home/christian97/.local/lib/python3.8/site-  
packages/sklearn/cross_decomposition/_pls.py:507: FutureWarning:
```

The attribute ``coef_`` will be transposed in version 1.3 to be consistent with other linear models in scikit-learn. Currently, ``coef_`` has a shape of `(n_features, n_targets)` and in the future it will have a shape of `(n_targets, n_features)`.

```
/home/christian97/.local/lib/python3.8/site-  
packages/sklearn/cross_decomposition/_pls.py:507: FutureWarning:
```

The attribute ``coef_`` will be transposed in version 1.3 to be consistent with other linear models in scikit-learn. Currently, ``coef_`` has a shape of `(n_features, n_targets)` and in the future it will have a shape of `(n_targets, n_features)`.

```
/home/christian97/.local/lib/python3.8/site-  
packages/sklearn/cross_decomposition/_pls.py:507: FutureWarning:
```

The attribute ``coef_`` will be transposed in version 1.3 to be consistent with other linear models in scikit-learn. Currently, ``coef_`` has a shape of `(n_features, n_targets)` and in the future it will have a shape of `(n_targets, n_features)`.

```
/home/christian97/.local/lib/python3.8/site-  
packages/sklearn/cross_decomposition/_pls.py:507: FutureWarning:
```

The attribute ``coef_`` will be transposed in version 1.3 to be consistent with other linear models in scikit-learn. Currently, ``coef_`` has a shape of `(n_features, n_targets)` and in the future it will have a shape of `(n_targets, n_features)`.

```
/home/christian97/.local/lib/python3.8/site-  
packages/sklearn/cross_decomposition/_pls.py:507: FutureWarning:
```

The attribute ``coef_`` will be transposed in version 1.3 to be consistent with other linear models in scikit-learn. Currently, ``coef_`` has a shape of `(n_features, n_targets)` and in the future it will have a shape of `(n_targets, n_features)`.

```
/home/christian97/.local/lib/python3.8/site-  
packages/sklearn/cross_decomposition/_pls.py:507: FutureWarning:
```

The attribute ``coef_`` will be transposed in version 1.3 to be consistent with other linear models in scikit-learn. Currently, ``coef_`` has a shape of `(n_features, n_targets)` and in the future it will have a shape of `(n_targets, n_features)`.

```
/home/christian97/.local/lib/python3.8/site-  
packages/sklearn/cross_decomposition/_pls.py:507: FutureWarning:
```

The attribute ``coef_`` will be transposed in version 1.3 to be consistent with other linear models in scikit-learn. Currently, ``coef_`` has a shape of `(n_features, n_targets)` and in the future it will have a shape of `(n_targets, n_features)`.

```
/home/christian97/.local/lib/python3.8/site-  
packages/sklearn/cross_decomposition/_pls.py:507: FutureWarning:
```

The attribute ``coef_`` will be transposed in version 1.3 to be consistent with other linear models in scikit-learn. Currently, ``coef_`` has a shape of `(n_features, n_targets)` and in the future it will have a shape of `(n_targets, n_features)`.

```
/home/christian97/.local/lib/python3.8/site-  
packages/sklearn/cross_decomposition/_pls.py:507: FutureWarning:
```

The attribute ``coef_`` will be transposed in version 1.3 to be consistent with other linear models in scikit-learn. Currently, ``coef_`` has a shape of `(n_features, n_targets)` and in the future it will have a shape of `(n_targets, n_features)`.

```
/home/christian97/.local/lib/python3.8/site-  
packages/sklearn/cross_decomposition/_pls.py:507: FutureWarning:
```

The attribute ``coef_`` will be transposed in version 1.3 to be consistent with other linear models in scikit-learn. Currently, ``coef_`` has a shape of `(n_features, n_targets)` and in the future it will have a shape of `(n_targets, n_features)`.

```
/home/christian97/.local/lib/python3.8/site-  
packages/sklearn/cross_decomposition/_pls.py:507: FutureWarning:
```

The attribute ``coef_`` will be transposed in version 1.3 to be consistent with other linear models in scikit-learn. Currently, ``coef_`` has a shape of `(n_features, n_targets)` and in the future it will have a shape of `(n_targets, n_features)`.

Visualization of VIP filtering

```
[102]: vip_metrics = make_subplots(rows=2, cols=3, shared_yaxes=True,
                                   subplot_titles=("Treshold 1: 1174 Features",
                                   ↪ "Treshold 1.5: 405 Features", "Treshold 2: 71 Features",
                                   ↪ "Annotated Features with Threshold",
                                   ↪ "Annotated Features with Threshold 1.5", "Annotated Features with",
                                   ↪ "Threshold 2"),
                                   vertical_spacing=0.09)

[103]: n_model = [1,2,3,4]

[104]: #####
#treshold 1
#####
vip_metrics.add_trace(go.Scatter(x=n_model, y=list(tr_1[0].values()),name="R2",
                                mode="lines+markers+text",
                                text=np.round(list(tr_1[0].values()),
                                ↪ decimals=3),
                                textposition="bottom center",
                                ↪ marker_color="#1f77b4"), row=1, col=1)
vip_metrics.add_trace(go.Scatter(x=n_model, y=list(tr_1[1].values()), name="Q2",
                                mode="lines+markers+text",
                                text=np.round(list(tr_1[1].values()),
                                ↪ decimals=3),
                                textposition="bottom center",
                                ↪ marker_color="#ff7f0e"), row=1, col=1)

vip_metrics.add_trace(go.Bar(x=["Total"], y=[70], name='Total',
                                ↪ marker_color="green"), row=2, col=1)
vip_metrics.add_trace(go.Bar(x=["HD"], y=[54], name='HD',
                                ↪ marker_color="purple"), row=2, col=1)
vip_metrics.add_trace(go.Bar(x=["KEGG"], y=[10], name='KEGG',
                                ↪ marker_color="orange"), row=2, col=1)
vip_metrics.add_trace(go.Bar(x=["LipidB"], y=[6], name='LipidBlast',
                                ↪ marker_color="blue"), row=2, col=1)

#####
#treshold 1.5
#####
vip_metrics.add_trace(go.Scatter(x=n_model, y=list(tr_1_5[0].values()),name="R2",
                                ↪ for treshold 1.5",
                                mode="lines+markers+text",
                                text=np.round(list(tr_1_5[0].values()),
                                ↪ decimals=3),
```

```

                                textposition="bottom center",
    ↪marker_color="#1f77b4", showlegend=False), row=1, col=2)
vip_metrics.add_trace(go.Scatter(x=n_model, y=list(tr_1_5[1].values()),name="Q2
    ↪for treshold 1.5",

                                mode="lines+markers+text",
                                text=np.round(list(tr_1_5[1].values()),

    ↪decimals=3),

                                textposition="bottom center",
    ↪marker_color="#ff7f0e", showlegend=False), row=1, col=2)

vip_metrics.add_trace(go.Bar(x=["Total"], y=[38], name='Total',
    ↪marker_color="green", showlegend=False), row=2, col=2)
vip_metrics.add_trace(go.Bar(x=["HD"], y=[31], name='HD',
    ↪marker_color="purple", showlegend=False), row=2, col=2)
vip_metrics.add_trace(go.Bar(x=["KEGG"], y=[7], name='KEGG',
    ↪marker_color="orange", showlegend=False), row=2, col=2)
vip_metrics.add_trace(go.Bar(x=["LipidB"], y=[0], name='LipidBlast',
    ↪marker_color="blue", showlegend=False), row=2, col=2)

#####
#treshold 2
#####
vip_metrics.add_trace(go.Scatter(x=n_model, y=list(tr_2[0].values()),name="R2
    ↪for treshold 2",

                                mode="lines+markers+text",
                                text=np.round(list(tr_2[0].values()),

    ↪decimals=3),

                                textposition="bottom center",
    ↪marker_color="#1f77b4", showlegend=False), row=1, col=3)
vip_metrics.add_trace(go.Scatter(x=n_model, y=list(tr_2[1].values()),name="Q2
    ↪for treshold 2",

                                mode="lines+markers+text",
                                text=np.round(list(tr_2[1].values()),

    ↪decimals=3),

                                textposition="bottom center",
    ↪marker_color="#ff7f0e", showlegend=False), row=1, col=3)

vip_metrics.add_trace(go.Bar(x=["Total"], y=[16], name='Total',
    ↪marker_color="green", showlegend=False), row=2, col=3)
vip_metrics.add_trace(go.Bar(x=["HD"], y=[12], name='HD',
    ↪marker_color="purple", showlegend=False), row=2, col=3)
vip_metrics.add_trace(go.Bar(x=["KEGG"], y=[4], name='KEGG',
    ↪marker_color="orange", showlegend=False), row=2, col=3)
vip_metrics.add_trace(go.Bar(x=["LipidB"], y=[0], name='LipidBlast',
    ↪marker_color="blue", showlegend=False), row=2, col=3)

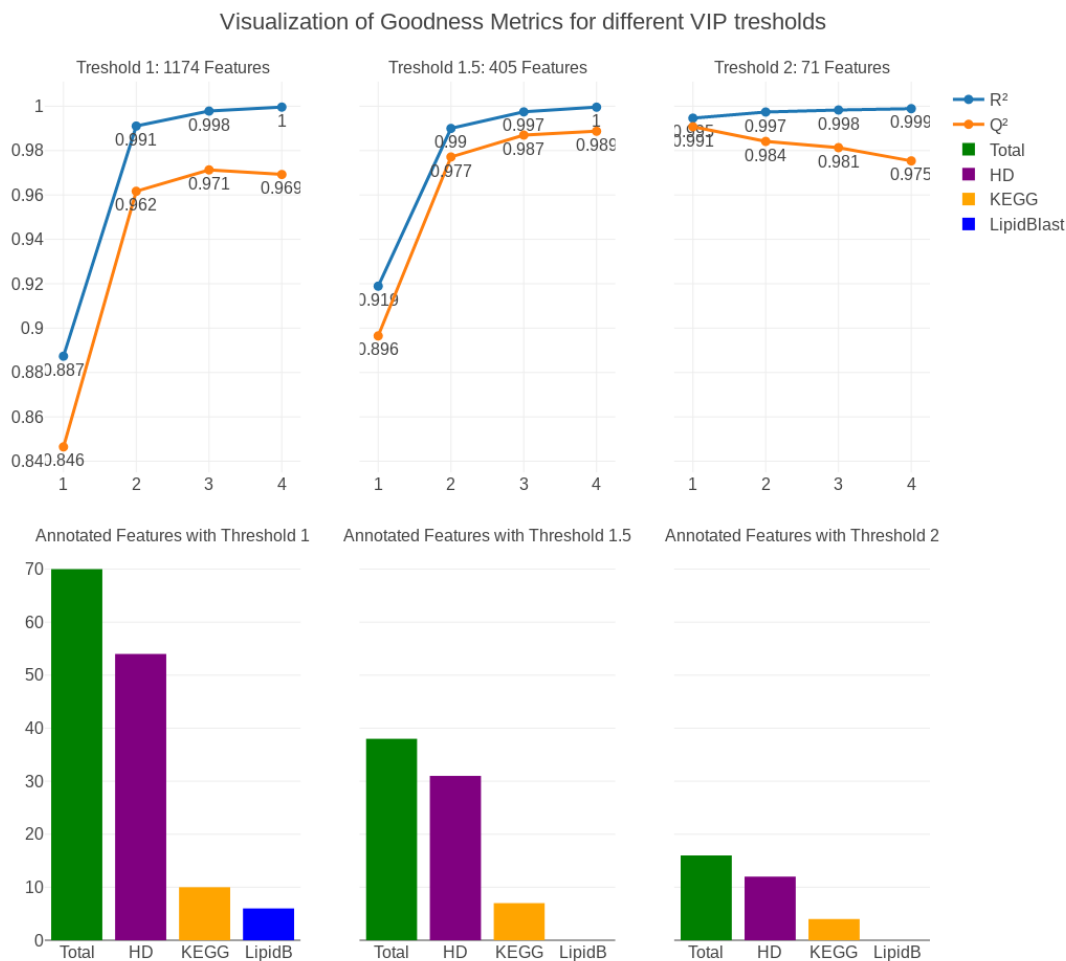
```



```

vip_metrics.update_layout(hovermode=False, font=dict(size=16), height=1000,
    width=1320,
    title_text="Visualization of Goodness Metrics for
different VIP thresholds",
    template="presentation")

```



The graphs at the top represent the goodness metrics for different $VIP_{thresholds}$, including the number of variables which were present in the model. The validation metrics of first model, with $VIP_{thresholds} = 1$, seemed to outperform the base model, with roughly one third of original 3250 variables. Consecutive models, reported validation metrics equally good; however, for $VIP_{thresholds} = 2$ the Q^2 metric started to decline when more than one LV was included in the model. Thus, the model obtained with a $VIP_{thresholds} = 2$ failed the model validation step.

At the bottom of the figure, it was represented the total number of annotated analytes which were present in the different models. Moreover, the graphs visualize the source of annotation, HD

referring to in-house DB, KEGG referring to Kyoto Encyclopedia of Genes and Genomes DB and LipidB referring to LipidBlast DB.

The number of annotated analytes obtained with $VIP_{thresholds} = 1$ was notably higher, 70, when compared to the number of annotated variables obtained with $VIP_{thresholds} = 1.5$, ~ 40 . Thus, the metabolic fingerprint extracted when setting a $VIP_{thresholds} = 1$ seemed to be more extensive, giving more chances of having highly informative digoxin-related biomarkers. For the final metabolic fingerprint only metabolites annotated at level 1, i.e HD annotation, were kept for graph-based analysis.

7 MetaboRank Output Analysis

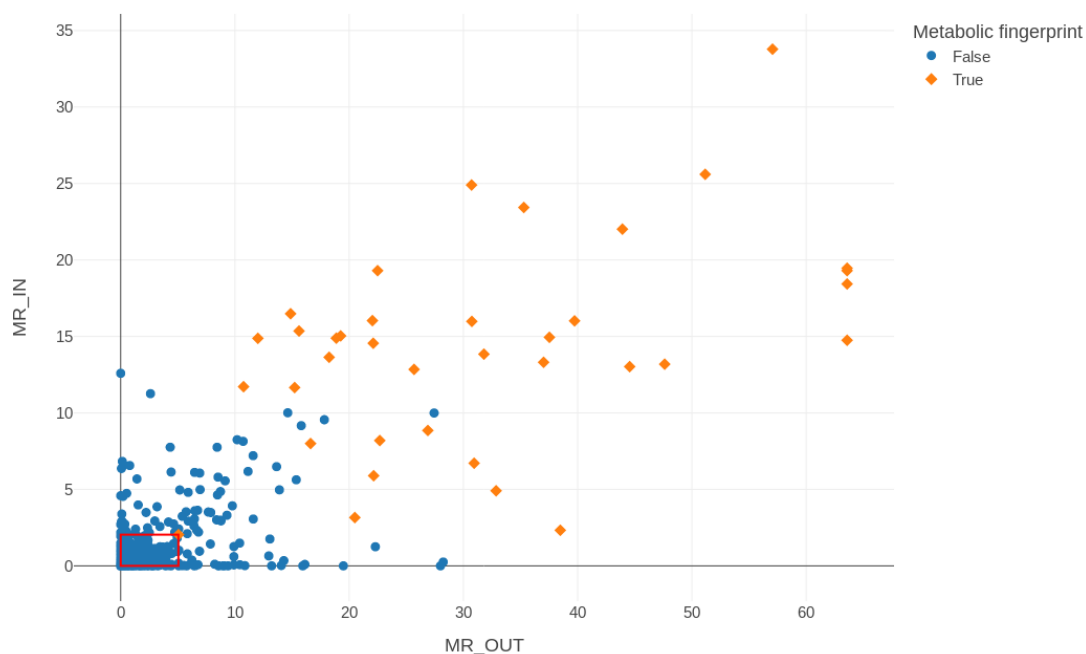
MetaboRank (MR) output consists of an excel file containing several spreadsheets which had to be filtered to collect data related to metabolites and probabilities computed by MR.

First, it was developed a function to open the excel file and retrieve a pandas dataframe. A list containing column names of interest is required.

```
[105]: def Read_MetRank_xls(path, sheet_name=1, col_names=None):  
        wk = open_workbook_xls(filename=path, ignore_workbook_corruption=True)  
        file = pd.read_excel(wk, sheet_name=sheet_name, names=col_names)  
        return file  
  
[106]: col_names=["Name", "Identifier", "Formula", "Monoisotopic_mass", "chebi",  
        ↪ "hmdb", "inchi", "inchikey", "kegg", "pubchem", "Fingerprint", "MR_OUT",  
        ↪ "MR_IN"]  
  
[107]: mr_output = Read_MetRank_xls("MR_noChemBackg.xls", col_names=col_names,  
        ↪ sheet_name=3)  
  
[108]: # filtering according to those entries which had different probability values  
        mr_output.query("MR_OUT != MR_IN", inplace=True)
```

A visualization similar to the one that can be found in MetExplore server is provided.

```
[109]: mr_plot=px.scatter(mr_output, x="MR_OUT", y="MR_IN",  
        ↪ color="Fingerprint", symbol="Fingerprint",  
        width=900, height=700, template="presentation",  
        ↪ hover_name="Identifier", hover_data=["Name"], symbol_sequence=["circle",  
        ↪ "diamond"])  
        mr_plot.update_layout(legend_title="Metabolic fingerprint", font=dict(size=15))  
        mr_plot.add_shape(type="rect", x0=0, x1=5.05, y0=0, y1=2.04,  
        ↪ line=dict(color="red", width=2))
```



To focus on a reduced subset of suggested candidates, the data point with lowest probabilities was taken as reference. It was located at the top right corner of the red squared. Hence, data points in the squared were filtered out.

```
[110]: subset = mr_output[(mr_output["MR_OUT"]>=5) | (mr_output["MR_IN"]>=2)]
subset =subset[subset["Fingerprint"]==False]
```

```
[111]: subset
```

```
[111]:
```

	Name	Identifier \
31	(2R,4S,5R,6R)-2-[(2R,3S,4R,5R)-5-(4-amino-2...	M_HC01162
36	(2S,3R)-2-azaniumyl-3-hydroxyoctadecyl phosphate	M_sph1p
41	(2S)-1-hydroxy-3-oxooctadecan-2-aminium	M_3dsphgn
55	(4-hydroxy-3-methoxyphenyl)acetaldehyde	M_3mox4hpac
59	(5-hydroxyindol-3-yl)acetaldehyde	M_5hoxindact
...
2556	Uracil	M_ura
2560	Uridine	M_uri
2582	Xanthine	M_xan
2583	xanthosine	M_xtsn
2584	Xanthosine 5-phosphate	M_xmp

	Formula	Monoisotopic_mass	chebi	hmdb \
31	C20H29N4O17P	628.127629	CHEBI:18098 CHEBI:58376	NaN

36	C18H39N05P	381.264410	CHEBI:16893 CHEBI:57939	HMDB01383
41	C18H38N02	299.282429	CHEBI:17862 CHEBI:58299	HMDB01480
55	C9H1003	166.062994	CHEBI:28111	HMDB05175
59	C10H9N02	175.063329	CHEBI:50157	HMDB04073
...
2556	C4H4N202	112.027277	CHEBI:17568	HMDB00300
2560	C9H12N206	244.069536	CHEBI:16704	HMDB00296
2582	C5H4N402	152.033425	CHEBI:17712 CHEBI:48517	HMDB00292
2583	C10H12N406	284.075684	CHEBI:18107	HMDB00299
2584	C10H11N409P	362.027462	CHEBI:15652 CHEBI:57464	HMDB01554

	inchi \
31	InChI=1S/C20H31N4017P/c21-10-1-2-24(19(35)22-1...
36	InChI=1S/C18H40N05P/c1-2-3-4-5-6-7-8-9-10-11-1...
41	InChI=1S/C18H37N02/c1-2-3-4-5-6-7-8-9-10-11-12...
55	InChI=1S/C9H1003/c1-12-9-6-7(4-5-10)2-3-8(9)11...
59	InChI=1S/C10H9N02/c12-4-3-7-6-11-10-2-1-8(13)5...
...	...
2556	InChI=1S/C4H4N202/c7-3-1-2-5-4(8)6-3/h1-2H,(H2...
2560	InChI=1S/C9H12N206/c12-3-4-6(14)7(15)8(17-4)11...
2582	InChI=1S/C5H4N402/c10-4-2-3(7-1-6-2)8-5(11)9-4...
2583	InChI=1S/C10H12N406/c15-1-3-5(16)6(17)9(20-3)1...
2584	InChI=1S/C10H13N409P/c15-5-3(1-22-24(19,20)21)...

	inchikey	kegg \
31	NaN	C03691
36	YHEDRJPUIRMZMP-ZWKOTPCHSA-M	C01120
41	KBUNOSOGGAARKZ-KRWDZBQOSA-O	C02934
55	GOQGGGANVKPMNH-UHFFFAOYSA-N	C05581
59	OBFAPCIUSYHFIE-UHFFFAOYSA-N	C05634
...
2556	ISAKRJJDGNUQOIC-UHFFFAOYSA-N	C00106 D00027 D09776
2560	DRTQHJPVMGBUCF-XVFCMESISA-N	C00299
2582	LRFVTTYWOQMYALW-UHFFFAOYSA-N	C00385
2583	UBORTCNDUKBEOP-UUOKFMHZSA-N	C01762
2584	NaN	C00655

	pubchem	Fingerprint \
31	656501	False
36	644260	False
41	439853	False
55	151276	False
59	74688	False
...
2556	10171240 10975456 11083870 11251987 ...	False
2560	1177 45356795 45358305 6029	False
2582	1188	False

2583	1189 45109822 64959	False
2584	73323	False

	MR_OUT	MR_IN
31	6.808849	2.207205
36	11.146692	6.172928
41	0.142130	6.827705
55	6.446680	0.000000
59	5.078950	0.000000
...
2556	16.120109	0.093043
2560	5.833207	0.783522
2582	8.728189	4.857052
2583	8.470786	4.635055
2584	4.634722	2.741113

[112 rows x 13 columns]

```
[112]: subset.shape
```

```
[112]: (112, 13)
```

In MR suggestion list can be included metabolites which are already annotated in the in-house DB. For this reason, a list of metabolites to be excluded was used for filtering. By this mean, it was ensured that the final subset only contained non-annotated metabolites.

```
[113]: chem_backg = pd.read_csv("MR_chem_background.txt", sep="\t",
    ↪names=["metabolite", "Identifier"])
```

```
[114]: candidates =subset[~subset["Identifier"].isin(chem_backg["Identifier"])]
```

```
[115]: candidates.shape
```

```
[115]: (79, 13)
```

Following step was to determine how many of the 79 candidates could be classified as **detectable**, **detected** or **non-annotated**

```
[116]: detectable_metabolites = subset[subset["Identifier"].
    ↪isin(chem_backg["Identifier"])]
```

Metabolites annotated in the experimental data were used to filter those metabolites already detected.

```
[117]: # this dataframe correspond to the MetExplore ID of annotated metabolites in
    ↪the in-house DB
experim_annot = pd.read_csv("xperim_annotation_info_ME-ID.txt",
    ↪names=["Identifier"])
```

```
[118]: detected_metabolites= subset[subset["Identifier"]
      ↪isin(experim_annot["Identifier"])]
```

There were 3 detectable metabolites, actually detected in the experimental data

- Homovanillate $VIP_{score} = 0.7960$
- 3-(4-hydroxyphenyl)pyruvate $VIP_{score} = 0.5927$
- methyl indole-3-acetate $VIP_{score} = 0.2097$

```
[119]: import plotly.graph_objects as go
```

```
[122]: labels=["non-annotated", "detectable", "detected"]
      values=[79, 33, 3]
```

```
[123]: donut = go.Figure(data=go.Pie(labels=labels, values=values, hole=.3))
      donut.update_traces(textinfo="value",
      ↪marker=dict(colors=['#1f77b4', '#ff7f0e', '#2ca02c']))
      donut.update_layout(width=800, height=500, title="Suggested candidates
      ↪classification", font=dict(size=16))
      #donut.write_image("../report_thesis/images/suggeted_classification.pdf")
      donut
```

Suggested candidates classification



After an analysis of MetaboRank output, it was retained a list of 79 non-annotated suggested candidates which were targeted for annotation. The list of non-annotated candidates corresponds to **candidates** dataframe.

Annotation at level 2 is usually done by comparing experimental MS/MS data with spectral databases. The benefit of HRMS lies on measurement of fragmentation pattern. Fragmentation pattern allows structure elucidation of the analytes. Given the fact that untargeted metabolomics studies detect tens of thousands of analytes, a focus was given on analytes which were selected for fragmentation, following standard DDA. Thus, annotation efforts were done solely on those analytes for which MS/MS data was available, aiming at achieving level 2 annotation. Selection of analytes to annotate was done by matching theoretical mass versus experimental mass. It was

used either molecular mass or m/z , depending on the kind of experimental mass available; always with a threshold of 5ppm in mass difference, To accurately perform mass matching it was essential to have correct molecular formula of neutral species. Hence, the list of 79 suggested candidates underwent thorough data curation for molecular formula, molecular mass as well as molecular descriptors, e.g. InChi, SMILES. Data curation was done by retrieving information from [PubChem DB](#), revealing that among the 79 candidates, 3 were actually proteins which were removed.

On one hand, mass matching was done using neutral molecular mass with those analytes which were reported by Progenesis software as having neutral mass; yielding no hit.

On the other hand, mass matching was done using m/z of analytes. For this approach, it was necessary to compute theoretical m/z for suggested candidates. Theoretical m/z computation was done using Molecular Spectrometry Adduct Calculator ([MSAC](#)) Python library. Adduct type selection was based on most common adduct in RPLC which included $[M+H]$, $[2M+H]$, $[M+Na]$, $[M+2H]$, $[M+3H]$ and $[M+H-H_2O]$.

7.1 Matching Neutral Masses

```
[124]: # filtering out proteins in MR output
candidates = candidates[candidates["Monoisotopic_mass"]<1000]

[125]: xperim_data = pd.read_csv("20201209_MGZ_720_RP.csv", sep="\t", header=4)

[128]: # getting the feature id which were annotated
annotated_features = xperim_data[~pd.isna(xperim_data["Accepted_ID"])]["Compound"]

[138]: def calculate_MassAccuracy(
    theoretical_df:pd.DataFrame, experimental_df:pd.
    DataFrame,colTheo_mass="Monoisotopic_mass", colExp_mass="Neutral mass (Da)",
    tolerance=5):
    """
    Description
    """
    if colExp_mass=="Neutral mass (Da)":
        # Getting the neutral mass for those features that where not annotated
        ("Accepted ID") and have MS fragmentation data
        expe_data = np.asarray(experimental_df[pd.
        isna(experimental_df["Accepted ID"])]~pd.isna(xperim_data["MSMS info_
        available"])]["Neutral mass (Da)"].dropna())
    else:
        expe_data = experimental_df[colExp_mass]
    def iterator(theor_vector):
        if expe_data[abs(1e6*(expe_data-theor_vector)/expe_data)<tolerance].
        any()==True:
            #obtaining the expected m/z value
            exp_hit= np.r_[expe_data[abs(1e6*(expe_data-theor_vector)/
            expe_data)<tolerance].values]
```

```

        #computing ppm mass error
        err_ppm = np.round(abs(1e6*(exp_hit-theor_vector)/exp_hit),
        ↪decimals=2)
        theo_hit = np.r_[theor_vector]
        # data manipulation
        if colTheo_mass== "Monoisotopic_mass":
            theo_filtered = theoretical_df[theoretical_df[colTheo_mass].
        ↪isin(theo_hit)].loc[:,["Name", "Formula", "Monoisotopic_mass"]].
        ↪sort_values(by="Monoisotopic_mass")
            exp_filtered = experimental_df[experimental_df[colExp_mass].
        ↪isin(exp_hit)].loc[:,["Compound", "Neutral mass (Da)"]].
        ↪sort_values(by="Neutral mass (Da)")
            final_df = pd.concat([theo_filtered.reset_index(drop=True),
        ↪exp_filtered.reset_index(drop=True)], axis="columns").loc[:,[
                "Name", "Formula", "Monoisotopic_mass", "Neutral mass
        ↪(Da)", "Compound"]]
        else:
            theo_filtered=theoretical_df[theoretical_df[colTheo_mass].
        ↪isin(theo_hit)].loc[:,[
                "Name", "adduct", "Monoisotopic_mass", "Formula", "adduct
        ↪mass", "smiles"]].sort_values(by="adduct mass")
            exp_filtered = experimental_df[experimental_df[colExp_mass].
        ↪isin(exp_hit)].loc[:, ["Compound", "m/z"]].sort_values(by="m/z")
            final_df = pd.concat([theo_filtered.reset_index(drop=True),
        ↪exp_filtered.reset_index(drop=True)], axis="columns").loc[:,[
                "Name", "Formula", "adduct mass", "m/z",
        ↪"adduct", "Compound", "smiles"]]
            final_df["error_ppm"]=err_ppm
            return final_df
    l = list(map(iterator, theoretical_df[colTheo_mass]))
    result = [i for i in l if i is not None]
    return result

```

```

[ ]: ppm_neutralMass =calculate_MassAccuracy(theoretical_df=candidates,
        ↪experimental_df=xperim_data,
                                                colTheo_mass="Monoisotopic_mass",
        ↪colExp_mass="Neutral mass (Da)")

```

```

[132]: len(ppm_neutralMass)

```

```

[132]: 0

```

There was no match with neutral masses

7.2 Matching Adduct Masses

```
[139]: #reading the adducts computed using MSAC library using the command line
adduct_list = pd.read_csv("msac_output.csv", sep=',')
```

```
[143]: # getting the feature id which were annotated
annotated_features = xperim_data[~pd.isna(xperim_data["Accepted_
↪ID"])]["Compound"]
```

```
[144]: ms_ms_spectra = pd.read_csv("MGZ_720_RP_listofMSMSforCP.csv", sep=",", header=1)
# filtering out features which were annotated
ms_ms_spectra= ms_ms_spectra[~ms_ms_spectra["Compound"].
↪isin(annotated_features)]
```

```
[145]: ppm_match =calculate_MassAccuracy(theoretical_df=adduct_list,
↪experimental_df=ms_ms_spectra, colTheo_mass="adduct mass", colExp_mass="m/z")
len(ppm_match)
```

```
[145]: 11
```

```
[146]: pd.concat(ppm_match, axis=0)
```

```
[146]:
```

		Name	Formula	adduct mass \
0		(2S)-1-hydroxy-3-oxooctadecan-2-aminium	C18H37NO2	322.271650
0		D-Mannose	C6H12O6	203.052609
0	1-(1,2,3,4,5-pentahydroxypent-1-yl)-1,2,3,4-te...	C17H22N2O7		123.054843
0	(5-hydroxyindol-3-yl)acetaldehyde	C10H9NO2		176.070605
0	1-(1,2,3,4,5-pentahydroxypent-1-yl)-1,2,3,4-te...	C17H22N2O7		367.149976
0	3,4-Dihydroxyphenylacetaldehyde	C8H8O3		153.054621
0	keto-phenylpyruvate	C9H8O3		165.054616
1		NaN	NaN	NaN
0		N-formyl-L-kynurenine	C11H12N2O4	237.086983
0	(2S)-1-hydroxy-3-oxooctadecan-2-aminium	C18H37NO2		282.279141
0	(4-hydroxy-3-methoxyphenyl)acetaldehyde	C9H10O3		149.059706
0	o-methylhippurate	C10H11NO3		176.070605

	m/z	adduct	Compound \
0	322.271423	M+Na	11.91_299.2823n
0	203.052204	M+Na	1.20_202.0449n
0	123.054950	M+3H	1.14_123.0549m/z
0	176.070160	M+H	7.33_176.0702m/z
0	367.149675	M+H	3.03_366.1424n
0	153.054111	M+H	5.36_170.0574n
0	165.054128	M+H	1.92_164.0471n
1	165.054160	NaN	1.14_164.0473n
0	237.086584	M+H	3.17_118.0397n
0	282.279024	M+H-H2O	14.95_281.2717n

```

0 149.059400 M+H-H2O 11.00_149.0594m/z
0 176.070160 M+H-H2O 7.33_176.0702m/z

```

	smiles	error_ppm
0	<chem>CCCCCCCCCCCCCCCC(=O)[C@H](CO)[NH2]</chem>	0.71
0	<chem>C([C@@H]1[C@H]([C@@H]([C@@H](C(O1)O)O)O)O)O</chem>	1.99
0	<chem>C1C(NC(C2=C1C3=CC=CC=C3N2)C(C(C(C(CO)O)O)O)O)C...</chem>	0.87
0	<chem>C1=CC2=C(C=C1O)C(=CN2)CC=O</chem>	2.53
0	<chem>C1C(NC(C2=C1C3=CC=CC=C3N2)C(C(C(C(CO)O)O)O)O)C...</chem>	0.82
0	<chem>C1=CC(=C(C=C1CC=O)O)O</chem>	3.33
0	<chem>C1=CC=C(C=C1)CC(=O)C(=O)O</chem>	2.96
1	NaN	2.77
0	<chem>C1=CC=C(C(=C1)C(=O)CC(C(=O)O)N)NC=O</chem>	1.68
0	<chem>CCCCCCCCCCCCCCCC(=O)[C@H](CO)[NH2]</chem>	0.41
0	<chem>COC1=C(C=CC(=C1)CC=O)O</chem>	2.05
0	<chem>CC1=CC=CC=C1C(=O)NCC(=O)O</chem>	2.53

Mass matching using m/z produced 9 unique hits, which were further examined. One out of the 9 hits, D-mannose, was checked again with the in-house DB, being classified as detectable but not detected. The reason for this incongruity was due to the time at which D-mannose was retained, $RT = 1.2$. In the context of LC, any analyte detected at $RT \leq 2$ meant that it was not retained by the column. Hence, it was considered that hits detected at $RT \leq 2$ were not suitable for annotation purposes as RPLC was not a suitable analytical platform for them.