

Example analysis of a publicly available dataset

This is an exploratory analysis of publicly available data from:

<https://www.metabolomicsworkbench.org/data/DRCCMetadata.php?Mode=Study&StudyID=ST000668&StudyType=MS&ResultType=1>
(<https://www.metabolomicsworkbench.org/data/DRCCMetadata.php?Mode=Study&StudyID=ST000668&StudyType=MS&ResultType=1>)

As requested, we do include this statement that comes with data:

This data is available at the NIH Common Fund's National Metabolomics Data Repository (NMDR) website, the Metabolomics Workbench, <https://www.metabolomicsworkbench.org> (<https://www.metabolomicsworkbench.org>), where it has been assigned Project ID PR000471. The data can be accessed directly via it's Project DOI: 10.21228/M85P6M This work is supported by NIH grant, U2C- DK119886.

Data preprocessing

In [1]:

```
1 # Author: stefano.manzini@gmail.com
2 # http://www.stemanz.altervista.org/
3 from matplotlib import pyplot as plt
4 %matplotlib inline
5 import seaborn as sns
6 from liputils import Lipid, make_residues_table
7 import pandas as pd
8 import numpy as np
9 from collections import Counter
10 import requests
```

Downloading the files from the repository in the local folder.

In [2]:

```
1 # data comes in two tables, from positive and negative ion modes
2
3 url1 = "https://www.metabolomicsworkbench.org/data/study_textformat_view.php?STU
4 url2 = "https://www.metabolomicsworkbench.org/data/study_textformat_view.php?STU
5
6 req1 = requests.get(url1, allow_redirects=False)
7 with open("ST000668_AN001022.csv", "wb") as o:
8     o.write(req1.content)
9
10 req2 = requests.get(url2, allow_redirects=False)
11 with open("ST000668_AN001023.csv", "wb") as o:
12     o.write(req2.content)
```

In [3]:

```
1 # manual inspection of the downloaded files is required to know how many rows
2 # need to be skipped to yield a non-malformed table
3 df1 = pd.read_csv("ST000668_AN001022.csv", sep="\t", skiprows=109)
4 df2 = pd.read_csv("ST000668_AN001023.csv", sep="\t", skiprows=109)
```

Even if it is hosted on Metabolomics Workbench, the is not 100% fully RefMet compliant. Before we use `liputils`, we want to use the online translator to convert all names into RefMet compliant ones. Before we do that, we must help the online translator by removing all unnecessary clutter from the analyte names.

In [4]:

```
1 # analytes contain information that prevent Metabolomics Workbench
2 # name translator from understanding who they are
3 # https://www.metabolomicsworkbench.org/databases/refmet/name_to_refmet.php
4 # example: CL 70:1; [M-2H](2- )@6.86
5 # it needs to be: CL 70:1
6
7 df1["Samples"] = df1["Samples"].apply(lambda x: x[:x.find(";")] if ";" in x else
8 df1["Samples"] = df1["Samples"].apply(lambda x: x[:x.find("(")] if "(" in x else
9 df2["Samples"] = df2["Samples"].apply(lambda x: x[:x.find(";")] if ";" in x else
10 df2["Samples"] = df2["Samples"].apply(lambda x: x[:x.find("(")] if "(" in x else
```

In [5]:

```
1 # the last two rows do not contain lipidomics data:
2 df1["Samples"][-3:]
```

Out[5]:

```
248 plasmenyl-PE 42:5
249 METABOLITES_END
250 #END
Name: Samples, dtype: object
```

In [6]:

```
1 df1 = df1.iloc[:-2,:]
2 df2 = df2.iloc[:-2,:]

```

In [7]:

```
1 #checking if columns are equal. If they are, we merge the two dataframes
2 for x,y in zip(df1.columns, df2.columns):
3     assert x == y

```

In [8]:

```
1 df = pd.concat([df1, df2], axis=0)
```

In [9]:

```
1 # this study uses a notation that is not RefMet compliant. We need to get
2 # online and retrieve translated terms from
3 # https://www.metabolomicsworkbench.org/databases/refmet/name_to_refmet.php

```

In [10]:

```
1 df = df.set_index("Samples")
```

In [11]:

```
1 # making a textfile with all analytes
2 with open("analytes.txt", "w") as o:
3     for analyte in df.index:
4         o.write(analyte+"\n")
```

In [12]:

```
1 # after manually inputting the analyte list in Metabolomics Workbench, a
2 # txt table of translated inputs is downloaded. We load this table:
3 translated = pd.read_csv("refmet_results.txt", sep="\t")
4 translated.head(10)
```

Out[12]:

	Input name	RefMet name	Formula	Mass	REGNO
0	Factors	NaN	NaN	NaN	NaN
1	CL 70:1	CL(70:1)	C79H152O17P2	1435.0505	NaN
2	CL 70:5	CL(70:5)	C79H144O17P2	1426.9879	NaN
3	CL 74:1	CL(74:1)	C83H160O17P2	1491.1131	NaN
4	CL 74:7	CL(74:7)	C83H148O17P2	1479.0192	NaN
5	CL 82:13	NaN	NaN	NaN	NaN
6	FFA(20:0)	Arachidic acid	C20H40O2	312.3028	41.0
7	FFA(20:1)	Heneicosenoic acid	C21H40O2	324.3028	NaN
8	FFA(20:2)	Eicosadienoic acid	C20H36O2	308.2715	NaN
9	FFA(22:0)	Behenic acid	C22H44O2	340.3341	43.0

In [13]:

```
1 # not all inputs were translated (Factors, on row 0, is an artifact from the old
2 translated.iloc[0,1] = "Factors" # replacing NaN with "Factors"
```

In [14]:

```
1 assert df.shape[0] == translated.shape[0]
```

In [15]:

```
1 # The tables have the same number of rows, so we can put the translated
2 # indices back into the former table, then drop all non translated
3 # compound (this is why we put back "Factors", as we still need to give labels
4 # to samples)
```

In [16]:

```
1 df["RefMet name"] = list(translated["RefMet name"])
2 df = df.dropna()
3 print(f"Total number of lipids: {df.shape[0]}")
```

Total number of lipids: 482

In [17]:

```
1 # it would be a good thing to rapidly identify each sample in terms of treatment
2 # as well as stratify data. Let's get to know the different experimental conditions
3 print(f"These are the possible experimental conditions:\n{'='*80}")
4 for c in set(df.loc["Factors",df.columns[:-1]]):
5     print(c)
```

These are the possible experimental conditions:

=====

Diet intervention:Baseline | Gender:male
Diet intervention:Unsaturated Fatty Acid | Gender:female
Diet intervention:Saturated Fatty Acid | Gender:female
Diet intervention:Baseline | Gender:female
Diet intervention:Overfeeding-2 weeks | Gender:female
Diet intervention:Overfeeding-2 weeks | Gender:male

In [18]:

```
1 # we are going to modify the table a bit, so we save some info about
2 # the samples for future reference
3 sample_book = {k:v for k,v in zip(df.columns[:-1], df.loc["Factors", df.columns[:-1]])}
```

In [19]:

```
1 # final preparation of the table for residue extraction. Setting the index with
2 # RefMet compliant lipid names
3 df = df.set_index("RefMet name")
```

In [20]:

```
1 df = df.drop(["Factors"])
```

Please note!

We need to **make sure that numeric data gets interpreted as numeric**, otherwise it will be discarded later on.

In [21]:

```
1 # liputils drops any non-numerical column by default.
2 # pandas has not interpreted the numbers in the table as numbers,
3 # so before proceeding we tell it that the table actually contains numbers.
4 df = df.astype("float64")
```

In [22]:

```
1 # This is what the final, polished table looks like
2 df.head()
```

Out[22]:

	S00019182	S00019183	S00019184	S00019185	S00019186	S00019187
RefMet name						
CL(70:1)	1.415543e+05	1.356717e+05	4.174906e+04	4.683934e+04	6.963792e+04	5.236800e+04
CL(70:5)	5.700773e+04	2.957225e+04	1.253442e+04	9.635603e+03	2.957225e+04	1.524695e+04
CL(74:1)	4.295867e+04	5.400946e+04	2.451510e+04	3.276228e+04	4.427461e+04	2.548709e+04
CL(74:7)	6.722344e+04	1.681784e+04	5.551705e+04	2.618583e+04	4.973006e+04	2.618583e+04
Arachidic acid	1.701166e+06	1.161941e+06	1.399456e+06	1.399456e+06	1.399456e+06	1.399456e+06

5 rows × 42 columns

Extraction of lipid residues with liputils

We directly feed our data table to `make_residues_table()`. This function returns a `pandas.DataFrame` and expects analytes as row index and samples as column index.

In [23]:

```
1 res = make_residues_table(df)
2
3 # residues are pulled in appearance order, so they need to be ordered
4 res = res.sort_index()
5 res.head()
```

Out[23]:

	S00019182	S00019183	S00019184	S00019185	S00019186	S00019187	S00019188
14:0	3.474830e+04	4.356682e+04	2.608958e+04	3.404043e+04	2.561295e+04	2.388370e+04	2.711111e+04
15:0	2.600530e+04	2.938390e+04	4.065718e+04	4.502925e+04	2.922490e+04	2.326481e+04	2.222222e+04
16:0	7.674743e+06	5.166182e+06	9.504088e+06	7.029468e+06	5.438244e+06	3.908616e+06	4.444444e+06
16:1	7.453597e+04	9.748588e+04	1.343240e+05	1.057079e+05	1.047244e+05	7.011630e+04	8.666667e+04
17:0	5.063393e+04	4.928738e+04	4.719015e+04	4.265817e+04	4.911188e+04	3.892206e+04	3.909091e+04

5 rows × 42 columns

In [24]:

```
1 print(f"A total of {len(res.index)} individual residues have been extracted.")
```

A total of 212 individual residues have been extracted.

In [25]:

```
1 # switching rows to columns for easier plotting
2 data = res.T
3
4 # giving some labels
5 def get_diet(stringlike):
6     if "Overfeeding-2 weeks" in stringlike:
7         return "Overfeeding"
8     elif "Baseline" in stringlike:
9         return "Baseline"
10    elif "Unsaturated" in stringlike:
11        return "Unsaturated fats"
12    elif "Saturated" in stringlike:
13        return "Saturated fats"
14    else:
15        return "unrecognized diet"
16
17
18 # adding some labels to data
19 data["sex"] = ["female" if "female" in sample_book.get(x) else "male" for x in data.index]
20 data["diet"] = [get_diet(sample_book.get(x)) for x in data.index]
```

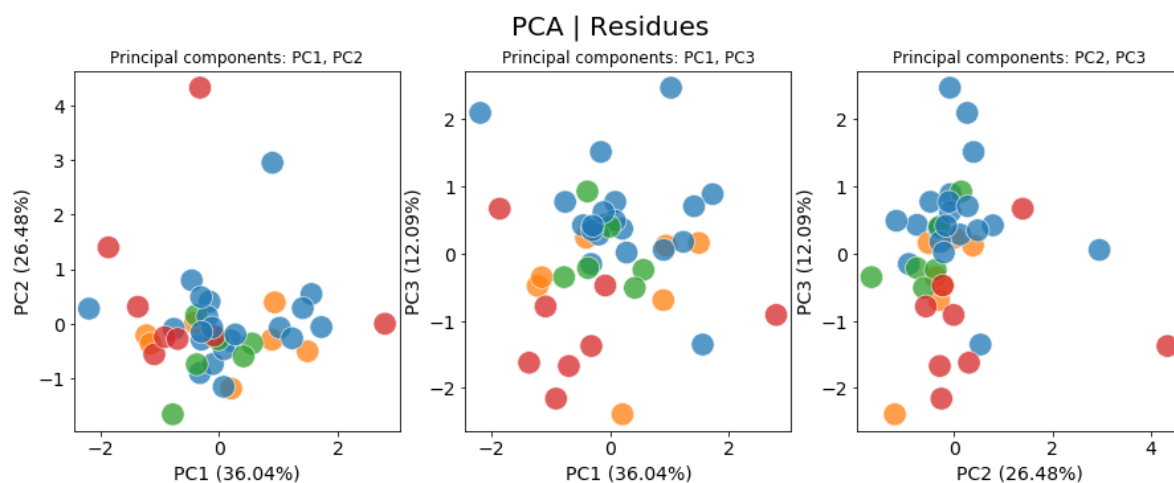
PCA | residues

In [26]:

```
1 from sklearn.decomposition import PCA
2
3 pca = PCA(3, whiten=True)
4
5 pc = pca.fit_transform(res.T)
6
7 repack = pd.DataFrame(pc, columns=["PC1", "PC2", "PC3"])
8 repack["samples"] = list(res.T.index)
9 repack["diet"] = [get_diet(sample_book.get(x)) for x in res.T.index]
10 repack["sex"] = ["female" if "female" in sample_book.get(x) else "male" for x in res.T.index]
11 repack = repack.set_index("samples")
```

In [27]:

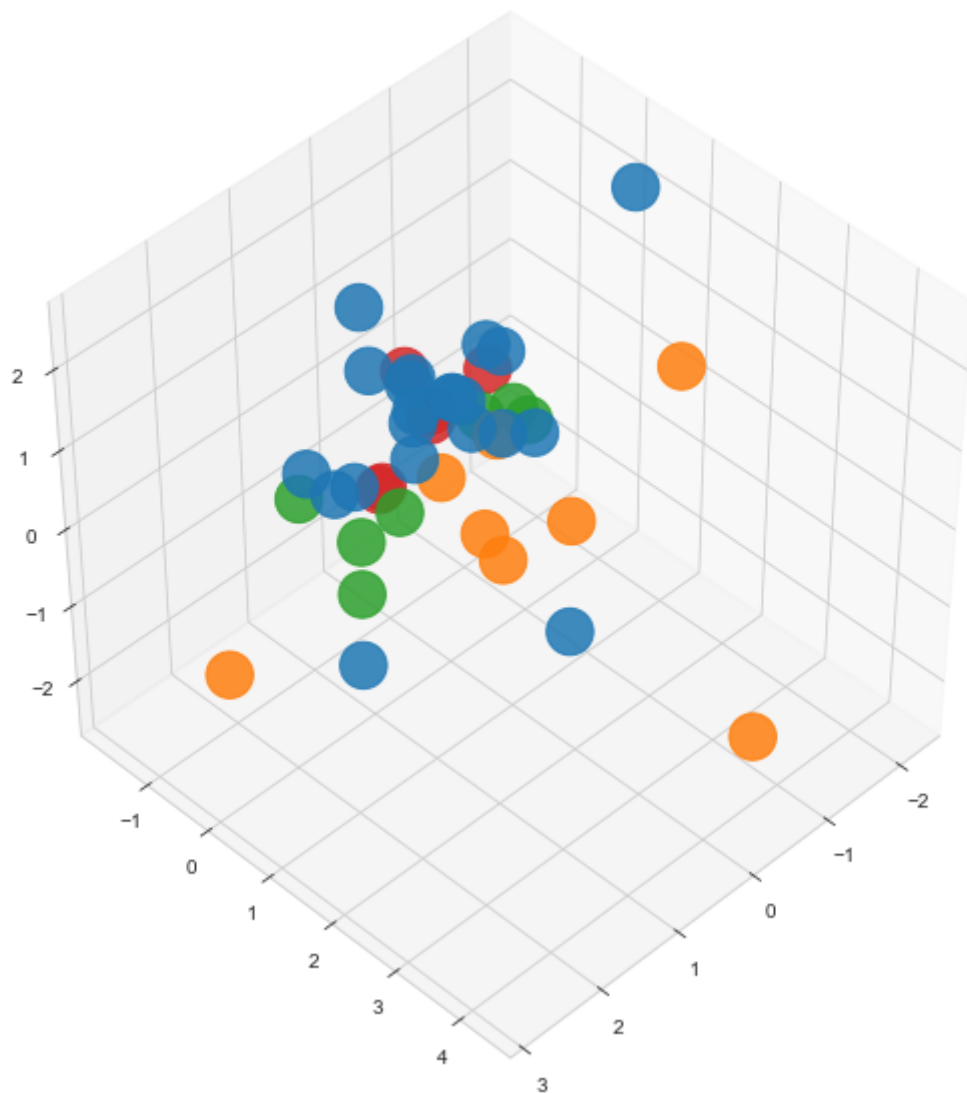
```
1 perc1, perc2, perc3 = pca.explained_variance_ratio_
2 expl_var = {"PC1": perc1,
3             "PC2": perc2,
4             "PC3": perc3,
5             }
6
7 # figure
8 fig, axes = plt.subplots(1, 3, figsize=(15,5))
9
10 for i, coord in enumerate(("PC1", "PC2"), ("PC1", "PC3"), ("PC2", "PC3")):
11     x, y = coord
12     axes[i].set_title(f"Principal components: {x}, {y}")
13
14     scattter = sns.scatterplot(
15         x=x, y=y, hue="diet", data=repack, s=300, ax=axes[i], alpha=.75
16     )
17
18     axes[i].set_xlabel(f"{x} ({round(expl_var.get(x)*100, 2)}%)", size=14)
19     axes[i].tick_params(axis='both', which='major', labelsize=14)
20     axes[i].set_ylabel(f"{y} ({round(expl_var.get(y)*100, 2)}%)", size=14)
21     axes[i].legend().set_visible(False)
22
23     fig.suptitle("PCA | Residues", size=20, y=1)
```



In [28]:

```
1 from mpl_toolkits.mplot3d import Axes3D
2
3 sns.set_style("whitegrid")
4
5 repacked = repack.copy() # saves time, code reuse
6 if True:
7     fig = plt.figure(figsize=(10,10))
8     ax = fig.add_subplot(111, projection='3d')
9
10    for g in set(repacked["diet"]):
11        #c = colors.get(g)
12        #m = marker.get(g)
13        #print(f"Group: {g}, color: {c}")
14        ax.scatter(
15            repacked[repacked["diet"] == g]["PC1"], # x
16            repacked[repacked["diet"] == g]["PC2"], # y
17            repacked[repacked["diet"] == g]["PC3"], # z
18            #c=c, #color
19            s=550, #size
20            alpha=.85,
21            #marker = m,
22        )
23
24    #ax.view_init(25, 33) #azim=35
25    ax.view_init(45, 45) #azim=35
26 print(f"Explained variance ratio: PC1: {round(perc1, 2)*100}%, PC2: {round(perc2, 2)*100}%, PC3: {round(perc3, 2)*100}%")
```

Explained variance ratio: PC1: 36.0%, PC2: 26.0%, PC3: 36.0%,



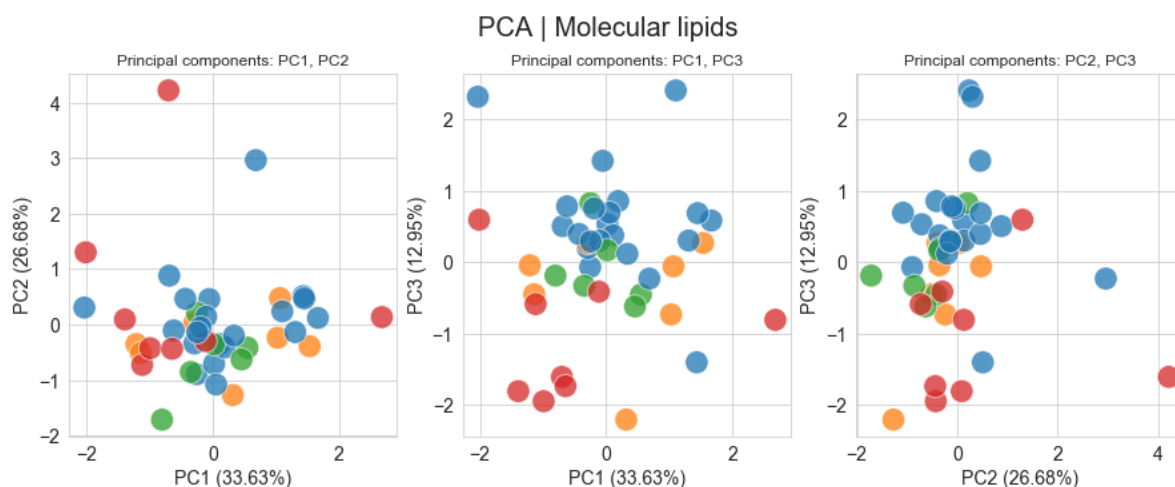
PCA | molecular lipids

In [29]:

```
1  pca = PCA(3, whiten=True)
2
3  pc = pca.fit_transform(df.T)
4
5  repack = pd.DataFrame(pc, columns=["PC1", "PC2", "PC3"])
6  repack["samples"] = list(df.T.index)
7  repack["diet"] = [get_diet(sample_book.get(x)) for x in df.T.index]
8  repack["sex"] = ["female" if "female" in sample_book.get(x) else "male" for x in df.T.index]
9  repack = repack.set_index("samples")
```

In [30]:

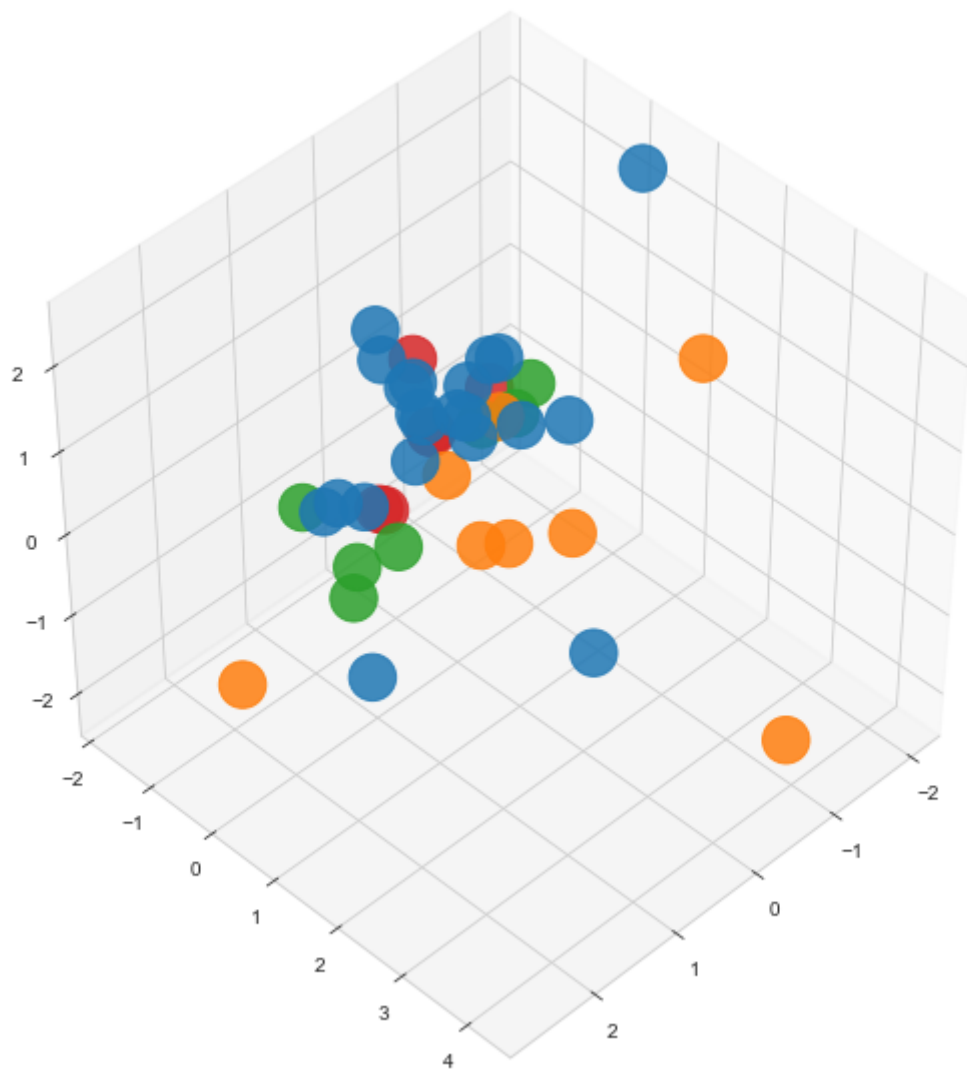
```
1 perc1, perc2, perc3 = pca.explained_variance_ratio_
2 expl_var = {"PC1": perc1,
3             "PC2": perc2,
4             "PC3": perc3,
5             }
6
7 # figure
8 fig, axes = plt.subplots(1, 3, figsize=(15,5))
9
10 for i, coord in enumerate(("PC1", "PC2"), ("PC1", "PC3"), ("PC2", "PC3")):
11     x, y = coord
12     axes[i].set_title(f"Principal components: {x}, {y}")
13
14     scattter = sns.scatterplot(
15         x=x, y=y, hue="diet", data=repack, s=300, ax=axes[i], alpha=.75
16     )
17
18     axes[i].set_xlabel(f"{x} ({round(expl_var.get(x)*100, 2)}%)", size=14)
19     axes[i].tick_params(axis='both', which='major', labelsize=14)
20     axes[i].set_ylabel(f"{y} ({round(expl_var.get(y)*100, 2)}%)", size=14)
21     axes[i].legend().set_visible(False)
22
23     fig.suptitle("PCA | Molecular lipids", size=20, y=1)
```



In [31]:

```
1 from mpl_toolkits.mplot3d import Axes3D
2
3 sns.set_style("whitegrid")
4
5
6 repacked = repack.copy() # saves time, code reuse
7 if True:
8     fig = plt.figure(figsize=(10,10))
9     ax = fig.add_subplot(111, projection='3d')
10
11     for g in set(repacked["diet"]):
12         #c = colors.get(g)
13         #m = marker.get(g)
14         #print(f"Group: {g}, color: {c}")
15         ax.scatter(
16             repacked[repacked["diet"] == g]["PC1"], # x
17             repacked[repacked["diet"] == g]["PC2"], # y
18             repacked[repacked["diet"] == g]["PC3"], # z
19             #c=c, #color
20             s=550, #size
21             alpha=.85,
22             #marker = m,
23         )
24     # curioso plot dei punti nella griglia..
25
26     #ax.view_init(25, 33) #azim=35
27     ax.view_init(45, 45) #azim=35
28 print(f"Explained variance ratio: PC1: {round(perc1, 2)*100}%, PC2: {round(perc2, 2)*100}%, PC3: {round(perc3, 2)*100}%")
```

Explained variance ratio: PC1: 34.0%, PC2: 27.0%, PC3: 33.6%,



Per-sample percentage change over baseline

A little more pre-processing

Each subject has two measures: baseline and treatment. We need to pair them to further the analysis. The files that are downloaded from the study contain this information, but it is cumbersome to read it with `pandas`. Instead we pull it from the [Project's page \(https://www.metabolomicsworkbench.org/data/subject_fetch.php?STUDY_ID=ST000668&STUDY_TYPE=MS&RESULT_TYPE=1\)](https://www.metabolomicsworkbench.org/data/subject_fetch.php?STUDY_ID=ST000668&STUDY_TYPE=MS&RESULT_TYPE=1):

There's no direct handle so we copypaste it into a new spreadsheet.

In [32]:

```
1 subjects = pd.read_csv("samples and subjects.txt", sep="\t")
```

In [33]:

```
1 subjects.head(2)
```

Out[33]:

	mb_sample_id	Subject name	Sample name	Diet intervention	Gender
0	SA038455	SU0019532	S00019182	Baseline	female
1	SA038467	SU0019533	S00019184	Baseline	female

Each sample has its own unique ID, but let's check whether there are two measurement per sample:

In [34]:

```
1 c = Counter(subjects["Subject name"])
2 c # for each subject, we have two measures
```

Out[34]:

```
Counter({'SU0019532': 2,
        'SU0019533': 2,
        'SU0019534': 2,
        'SU0019535': 2,
        'SU0019536': 2,
        'SU0019537': 2,
        'SU0019538': 2,
        'SU0019539': 2,
        'SU0019540': 2,
        'SU0019541': 2,
        'SU0019542': 2,
        'SU0019543': 2,
        'SU0019544': 2,
        'SU0019545': 2,
        'SU0019548': 2,
        'SU0019549': 2,
        'SU0019553': 2,
        'SU0019546': 2,
        'SU0019547': 2,
        'SU0019550': 2,
        'SU0019551': 2,
        'SU0019552': 2})
```

We now need to process our data so that the columns that hold the samples are paired, then the percentage change is calculated.

In [35]:

```
1 subjects = subjects.set_index("Subject name")
2
3 subjects_book = {}
4 for k in c:
5     subjects_book.setdefault(k, list(subjects.loc[k]["Sample name"]))
```

In [36]:

```
1 subjects_book
```

Out[36]:

```
{'SU0019532': ['S00019182', 'S00019183'],
'SU0019533': ['S00019184', 'S00019185'],
'SU0019534': ['S00019186', 'S00019187'],
'SU0019535': ['S00019188', 'S00019189'],
'SU0019536': ['S00019190', 'S00019191'],
'SU0019537': ['S00019192', 'S00019193'],
'SU0019538': ['S00019194', 'S00019195'],
'SU0019539': ['S00019196', 'S00019197'],
'SU0019540': ['S00019198', 'S00019199'],
'SU0019541': ['S00019200', 'S00019201'],
'SU0019542': ['S00019202', 'S00019203'],
'SU0019543': ['S00019204', 'S00019205'],
'SU0019544': ['S00019206', 'S00019207'],
'SU0019545': ['S00019208', 'S00019209'],
'SU0019548': ['S00019214', 'S00019215'],
'SU0019549': ['S00019216', 'S00019217'],
'SU0019553': ['S00019224', 'S00019225'],
'SU0019546': ['S00019210', 'S00019211'],
'SU0019547': ['S00019212', 'S00019213'],
'SU0019550': ['S00019218', 'S00019219'],
'SU0019551': ['S00019220', 'S00019221'],
'SU0019552': ['S00019222', 'S00019223']}
```

Smartly enough, the baseline conditions are even numbered, and post-treatment values are odd numbered. So we can easily check above that the samples were correctly paired.

In [37]:

```
1 templist = []
2 for i, k in enumerate(subjects_book.keys()):
3     try:
4         col1, col2 = subjects_book.get(k)
5         base, treatment = res[col1], res[col2]
6         templist.append(((treatment - base)/base)*100)
7     except KeyError:
8         print(f"Warning: {k} can't be found among samples ?!")
9 print(f"Processed {i+1} entries")
```

```
Warning: SU0019540 can't be found among samples ?!
Processed 22 entries
```

One pair of samples is listed in the project's files, but it is not present in the project's data.

Final prepping of data before plotting:

In [38]:

```
1 residue_perc = pd.concat(templist, axis=1)
2 newcols = {old: new for old, new in zip(residue_perc.columns, subjects_book.keys)}
3
4 residue_perc = residue_perc.rename(columns=newcols)
5
6 dfperc = residue_perc.T
7 dfperc["diet"] = [list(subjects.loc[key]["Diet intervention"])[1] for key in dfperc.keys()]
```

In [39]:

```
1 paper_names = {
2     "Saturated Fatty Acid" : "Saturated fats",
3     "Unsaturated Fatty Acid" : "Unsaturated fats",
4     "Overfeeding-2 weeks" : "Overfeeding",
5 }
```

In [40]:

```
1 dfperc["diet"] = dfperc["diet"].map(paper_names)
```

Percent change over baseline plots

We are using two `liputils` functions that help managing the residues that we want to plot:

`saturated()` returns `True` if a residue is saturated, and `False` otherwise. This is convenient as unsaturated residues can be specified via the `not` operator, as in `not saturated(residue)`.

`max_carbon()` returns `True` if the specified residue is up to the specified number of carbon atoms in length, `False` otherwise.

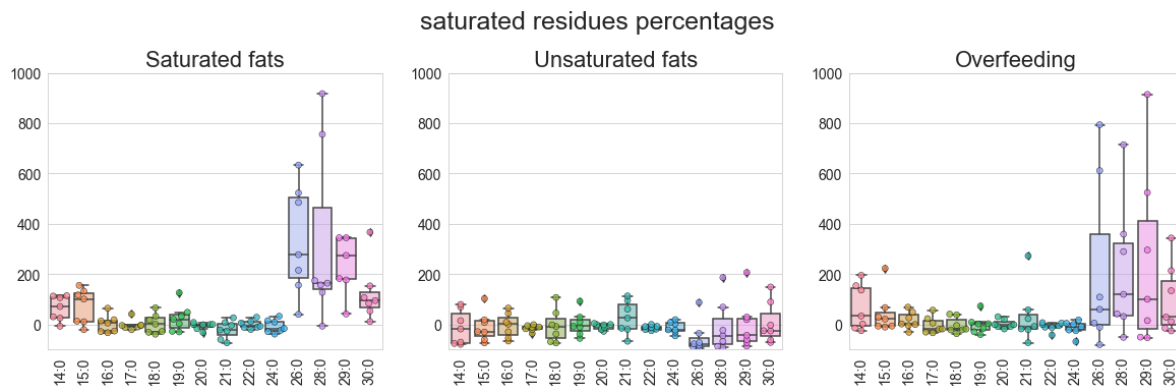
They can be combined to pythonically define the residue subset that needs plotting, like in here (line 22):

In [41]:

```
1 from liputils import saturated, max_carbon
```

In [42]:

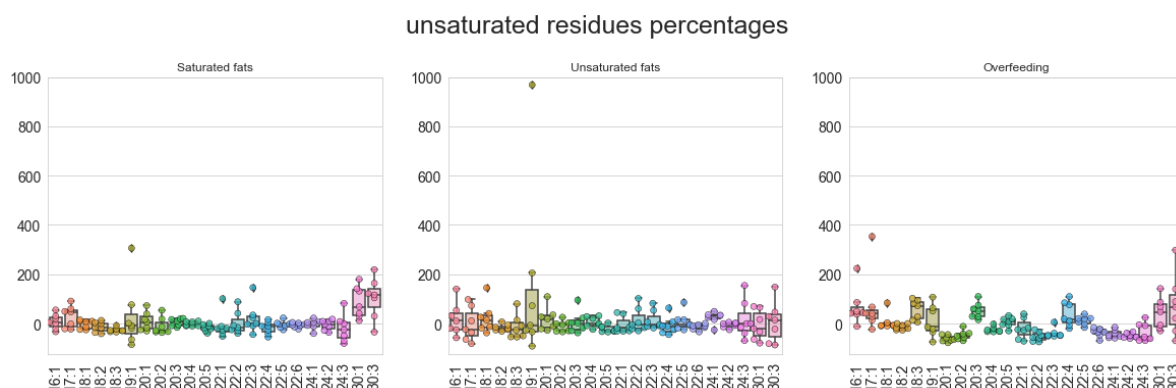
```
1 case = "saturated residues percentages"
2
3 fig, axes = plt.subplots(1, 3, figsize=(20,5))
4 label = "diet"
5 labelled_columns = 1
6 TITLE_SIZE = 22
7 TICKSIZE = 14
8
9 plt.suptitle(case, y=1.05, size=TITLE_SIZE*1.13)
10 for i, group in enumerate(('Saturated fats', 'Unsaturated fats', 'Overfeeding')):
11
12     axes[i].set_title(f"{group}", size=TITLE_SIZE)
13
14     # applying a threshold.
15     # =====
16     cols = abs(dfperc._get_numeric_data().max()) >= 0 # that is, all of them
17     cols_to_keep = cols[cols].index
18     table = dfperc[cols_to_keep]
19
20     # saturated or insaturated?
21     # =====
22     table = table.loc[:, [saturated(x) and max_carbon(x, 30) for x in table.columns]]
23
24     # putting back sample labels
25     table = pd.concat([table, dfperc[dfperc.columns[-labelled_columns:]]], axis=1)
26
27     table=table[table[label] == group]
28
29     # we needed to put back the labels to drop non-group rows, now
30     # we just need to plot the numbers
31     table=table[table.columns[:-labelled_columns]] # only residue columns
32
33     #bar = sns.barplot(data=table, ax=axes[i])
34     box = sns.boxplot(data=table, ax=axes[i], whis=1.5)
35     for patch in box.artists:
36         r, g, b, a = patch.get_facecolor()
37         patch.set_facecolor((r, g, b, .5))
38     swarm = sns.swarmplot(data=table, ax=axes[i], s=6, alpha=.8, edgecolor="black")
39     axes[i].tick_params(rotation=90, labelsz=TICKSIZE, axis="x")
40     axes[i].tick_params(labelz=TICKSIZE, axis="y")
41     axes[i].set_ylim(-100, 1000)
```



In [43]:

```
1 case = "unsaturated residues percentages"
2
3 fig, axes = plt.subplots(1, 3, figsize=(20,5))
4 label = "diet"
5 labelled_columns = 1
6 TITLE_SIZE = 22
7 TICKSIZE = 14
8
9 plt.suptitle(case, y=1.05, size=TITLE_SIZE*1.13)
10 for i, group in enumerate(('Saturated fats', 'Unsaturated fats', 'Overfeeding')):
11
12     axes[i].set_title(f"{group}")
13
14     # applying a threshold.
15     # =====
16     cols = abs(dfperc._get_numeric_data().max()) >= 0 # that is, all of them
17     cols_to_keep = cols[cols].index
18     table = dfperc[cols_to_keep]
19
20     # saturated or insaturated?
21     # =====
22     table = table.loc[:, [not saturated(x) and max_carbon(x, 30) for x in table.columns]]
23
24     # putting back sample labels
25     table = pd.concat([table, dfperc[dfperc.columns[-labelled_columns:]]], axis=1)
26
27     table=table[table[label] == group]
28
29     # we needed to put back the labels to drop non-group rows, now
30     # we just need to plot the numbers
31     table=table[table.columns[:-labelled_columns]] # only residue columns
32
33     #bar = sns.barplot(data=table, ax=axes[i])
34     box = sns.boxplot(data=table, ax=axes[i], whis=2)
35     for patch in box.artists:
36         r, g, b, a = patch.get_facecolor()
37         patch.set_facecolor((r, g, b, .5))
38     swarm = sns.swarmplot(data=table, ax=axes[i], s=6, alpha=.8, edgecolor="black")
39     axes[i].tick_params(rotation=90, labelsz=TICKSIZE, axis="x")
40     axes[i].tick_params(labelsize=TICKSIZE, axis="y")
41     axes[i].set_ylim(-125, 1000)
42
43 print(case)
44 #plt.savefig(f"/Users/manz/Desktop/{case}.png", dpi=600, bbox_inches="tight")
```

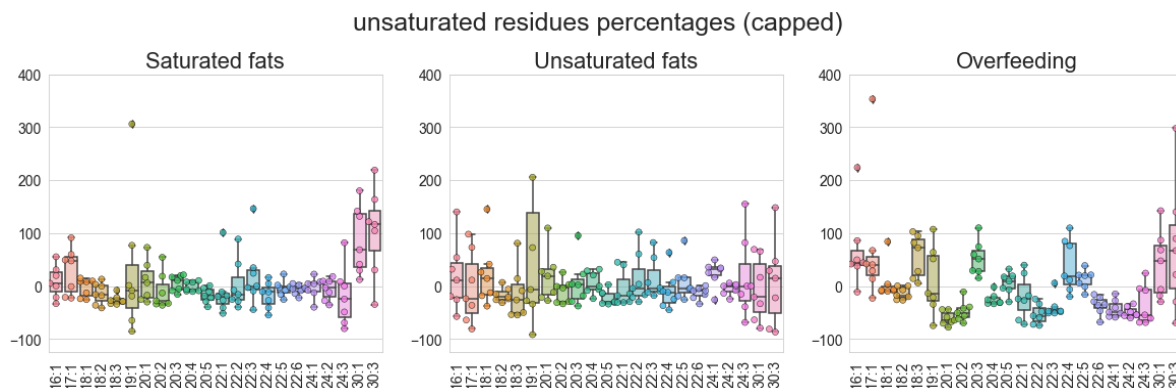
unsaturated residues percentages



There's a clear outlier in the 19:1 residue in Unsaturated fats, that renders the figure hard to read, flattening all the boxes. Let's reduce the span of y-axis:

In [44]:

```
1 case = "unsaturated residues percentages (capped)"
2
3 fig, axes = plt.subplots(1, 3, figsize=(20,5))
4 label = "diet"
5 labelled_columns = 1
6 TITLE_SIZE = 22
7 TICKSIZE = 14
8
9 plt.suptitle(case, y=1.05, size=TITLE_SIZE*1.13)
10 for i, group in enumerate(('Saturated fats', 'Unsaturated fats', 'Overfeeding')):
11
12     axes[i].set_title(f"{group}", size=TITLE_SIZE)
13
14     # applying a threshold.
15     # =====
16     cols = abs(dfperc._get_numeric_data().max()) >= 0 # that is, all of them
17     cols_to_keep = cols[cols].index
18     table = dfperc[cols_to_keep]
19
20     # saturated or insaturated?
21     # =====
22     table = table.loc[:, [not saturated(x) and max_carbon(x, 30) for x in table.columns]]
23
24     # putting back sample labels
25     table = pd.concat([table, dfperc[dfperc.columns[-labelled_columns:]]], axis=1)
26
27     table=table[table[label] == group]
28
29     # we needed to put back the labels to drop non-group rows, now
30     # we just need to plot the numbers
31     table=table[table.columns[:-labelled_columns]] # only residue columns
32
33     #bar = sns.barplot(data=table, ax=axes[i])
34     box = sns.boxplot(data=table, ax=axes[i], whis=2)
35     for patch in box.artists:
36         r, g, b, a = patch.get_facecolor()
37         patch.set_facecolor((r, g, b, .5))
38     swarm = sns.swarmplot(data=table, ax=axes[i], s=6, alpha=.8, edgecolor="black")
39     axes[i].tick_params(rotation=90, labelsz=TICKSIZE, axis="x")
40     axes[i].tick_params(labelz=TICKSIZE, axis="y")
41     axes[i].set_ylim(-125, 400)
```



This plot look definitely better than the previous one.

Residues bar plots

Baseline and treatment

Residue abundance vary greatly. An arbitrary threshold has been chosen to differentiate "abundant" and "less abundant" residues. Charts thus have different y-axis maxima, just for the sake of clarity.

Plotting saturated residues

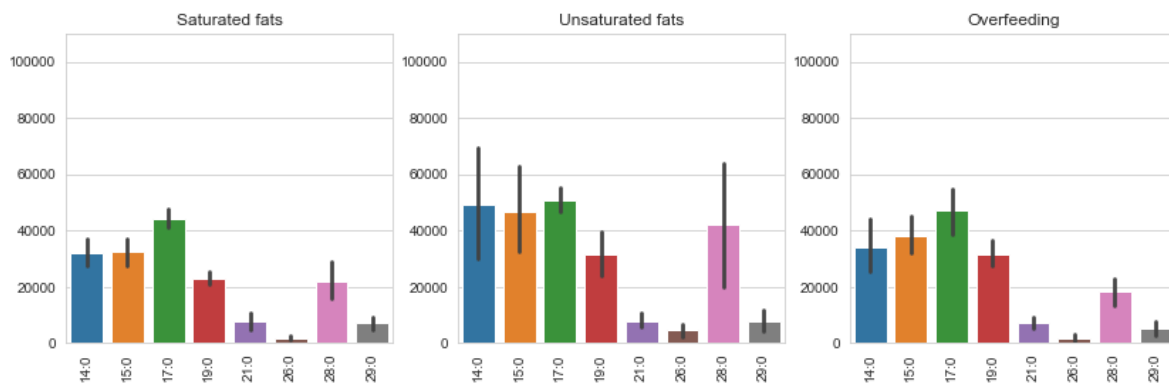
In [45]:

```
1 # taken from online sample details
2 baseline_samples = {
3     "Overfeeding" : ['S00019214', 'S00019216', 'S00019218', 'S00019220',
4                     'S00019222', 'S00019224'],
5     "Unsaturated fats" : ['S00019196', 'S00019200', 'S00019202', 'S00019204',
6                          'S00019206', 'S00019208'],
7     "Saturated fats" : ['S00019182', 'S00019184', 'S00019186', 'S00019188',
8                        'S00019190', 'S00019192', 'S00019194'],
9 }
```

In [46]:

```
1 label = "Baseline conditions, abundant, saturated residues"
2
3 fig, axes = plt.subplots(1, 3, figsize=(14,4))
4 plt.suptitle(label, size=18, y=1.05)
5 labelled_columns = 2
6
7 for i, key in enumerate(("Saturated fats", "Unsaturated fats", "Overfeeding")):
8
9     axes[i].set_title(f"{key}")
10
11     table = data.loc[baseline_samples.get(key), :].copy()
12
13     # saturated or insaturated?
14     # =====
15     table = table.loc[:, [saturated(x) and max_carbon(x, 30) for x in table._get_
16
17     cols = table.max() < 100_000
18     cols = cols[cols].index
19     table = table[cols]
20     BASELINE_COLS = cols.copy()
21
22     bar = sns.barplot(data=table, ax=axes[i],
23                       #alpha=.5
24                       )
25     #swarm = sns.swarmplot(data=table, ax=axes[i], s=10, alpha=.85, edgecolor="k")
26     axes[i].tick_params(rotation=90, labelsz=10, axis="x")
27     axes[i].tick_params(labelsz=10, axis="y")
28     axes[i].set_ylim(0, 110_000)
```

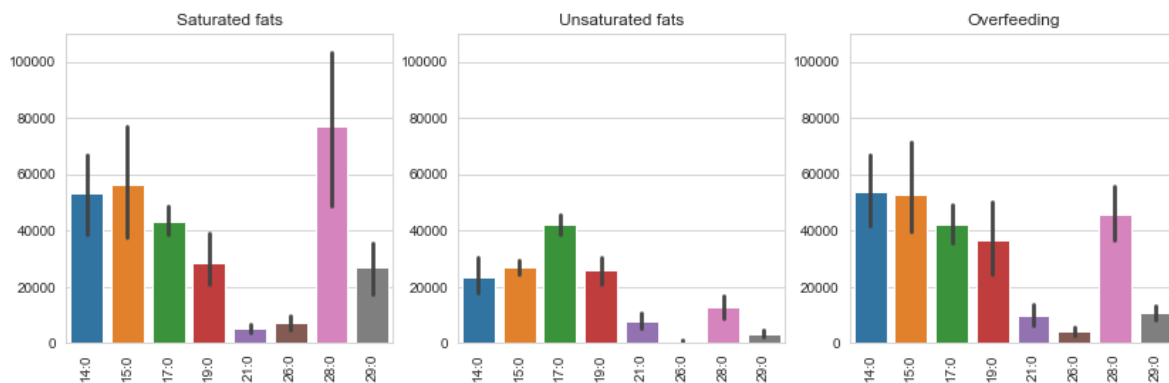
Baseline conditions, abundant, saturated residues



In [47]:

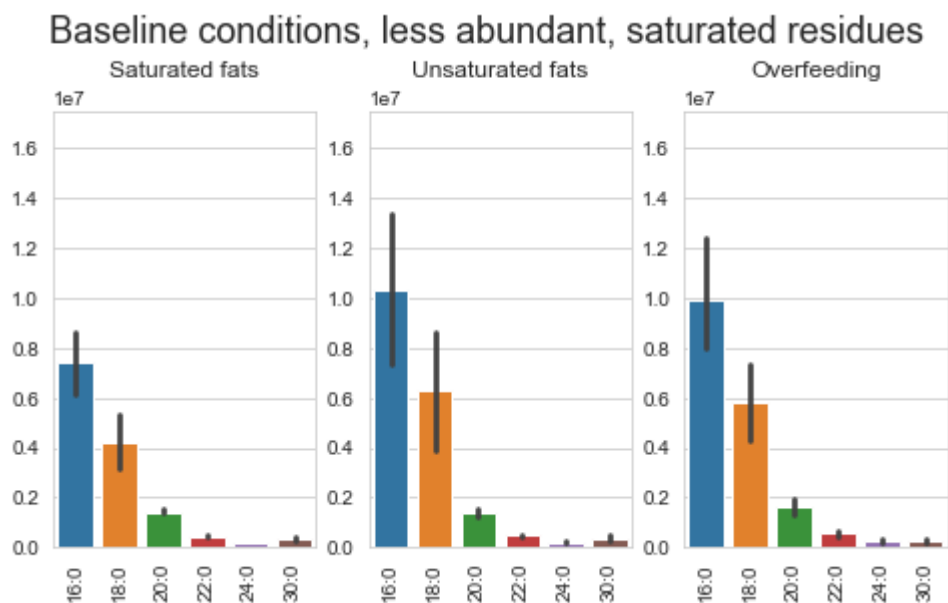
```
1 label = "After treatment, abundant, saturated residues"
2
3 fig, axes = plt.subplots(1, 3, figsize=(14,4))
4 plt.suptitle(label, size=18, y=1.05)
5 labelled_columns = 2
6
7 for i, key in enumerate(("Saturated fats", "Unsaturated fats", "Overfeeding")):
8
9     axes[i].set_title(f"{key}")
10
11     #table = data.loc[baseline_samples.get(key), :].copy()
12     table = data[data["diet"] == key].copy()
13
14     # saturated or insaturated?
15     # =====
16     table = table.loc[:, [saturated(x) and max_carbon(x, 30) for x in table._get_
17
18     # the same as previous plot
19     table = table[BASLINE_COLS]
20
21     bar = sns.barplot(data=table, ax=axes[i],
22                       #alpha=.5
23                       )
24     #swarm = sns.swarmplot(data=table, ax=axes[i], s=10, alpha=.85, edgecolor="k")
25     axes[i].tick_params(rotation=90, labelsize=10, axis="x")
26     axes[i].tick_params(labelsize=10, axis="y")
27     axes[i].set_ylim(0, 110_000)
```

After treatment, abundant, saturated residues



In [48]:

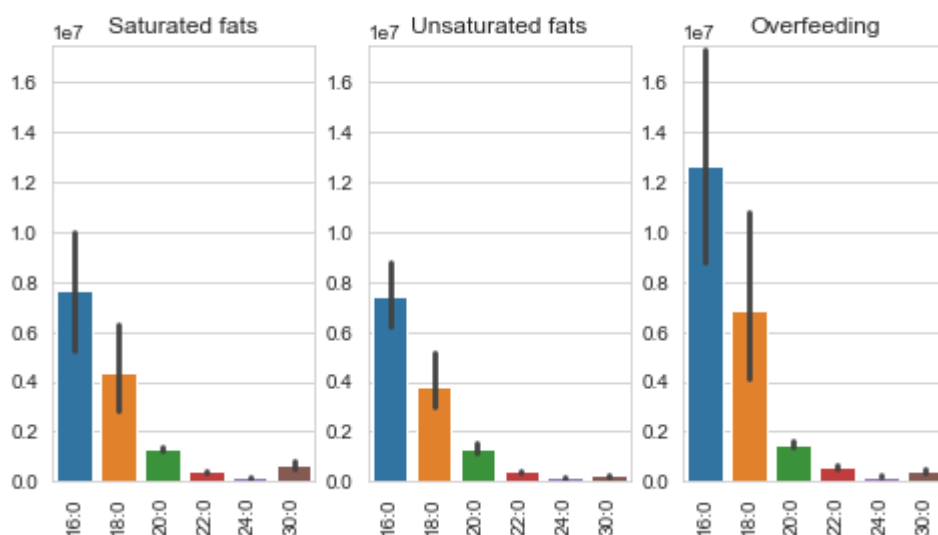
```
1 label = "Baseline conditions, less abundant, saturated residues"
2
3 fig, axes = plt.subplots(1, 3, figsize=(8,4))
4 plt.suptitle(label, size=18, y=1.05)
5 labelled_columns = 2
6
7 for i, key in enumerate(("Saturated fats", "Unsaturated fats", "Overfeeding")):
8
9     axes[i].set_title(f"{key}", y=1.05)
10
11     table = data.loc[baseline_samples.get(key), :].copy()
12
13     # saturated or unsaturated?
14     # =====
15     table = table.loc[:, [saturated(x) and max_carbon(x, 30) for x in table._get_
16
17     cols = table.max() >= 100_000
18     cols = cols[cols].index
19     table = table[cols]
20     BASELINE_COLS = cols.copy()
21
22     bar = sns.barplot(data=table, ax=axes[i],
23                       #alpha=.5
24                       )
25     #swarm = sns.swarmplot(data=table, ax=axes[i], s=10, alpha=.85, edgecolor="k")
26     axes[i].tick_params(rotation=90, labelsize=10, axis="x")
27     axes[i].tick_params(labelsize=10, axis="y")
28     axes[i].set_ylim(0, 17_500_000)
```



In [49]:

```
1 label = "After treatment, less abundant, saturated residues"
2
3 fig, axes = plt.subplots(1, 3, figsize=(8,4))
4 plt.suptitle(label, size=18, y=1.05)
5 labelled_columns = 2
6
7 for i, key in enumerate(("Saturated fats", "Unsaturated fats", "Overfeeding")):
8
9     axes[i].set_title(f"{key}")
10
11     #table = data.loc[baseline_samples.get(key), :].copy()
12     table = data[data["diet"] == key].copy()
13
14     # saturated or insaturated?
15     # =====
16     table = table.loc[:, [saturated(x) and max_carbon(x, 30) for x in table._get_
17
18     # the same as previous plot
19     table = table[BASLINE_COLS]
20
21     bar = sns.barplot(data=table, ax=axes[i],
22                       #alpha=.5
23                       )
24     #swarm = sns.swarmplot(data=table, ax=axes[i], s=10, alpha=.85, edgecolor="k")
25     axes[i].tick_params(rotation=90, labels=10, axis="x")
26     axes[i].tick_params(labels=10, axis="y")
27     axes[i].set_ylim(0, 17_500_000)
```

After treatment, less abundant, saturated residues

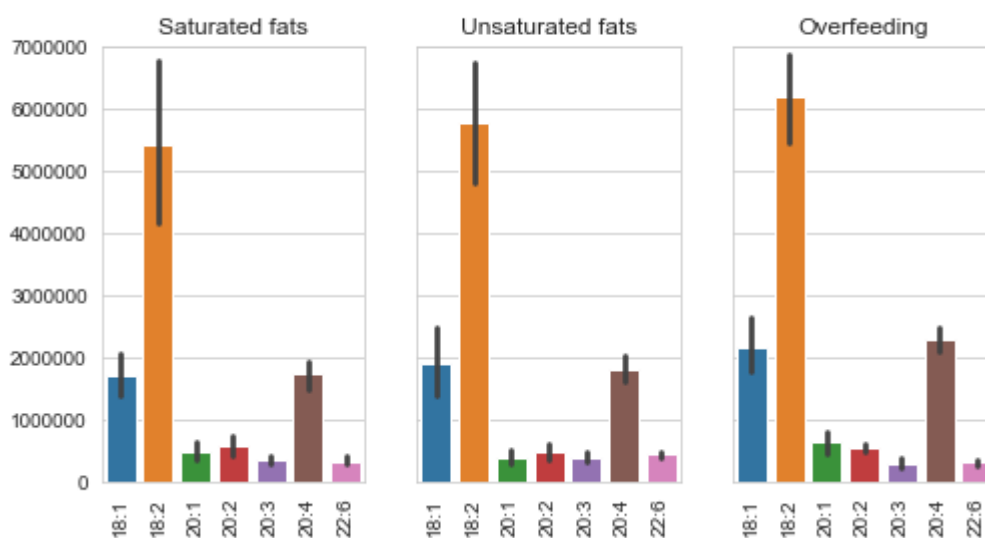


Plotting unsaturated residues

In [50]:

```
1 label = "Baseline conditions, abundant, unsaturated residues"
2
3 fig, axes = plt.subplots(1, 3, figsize=(8,4))
4 plt.suptitle(label, size=18, y=1.05)
5 labelled_columns = 2
6
7 for i, key in enumerate(("Saturated fats", "Unsaturated fats", "Overfeeding")):
8
9     axes[i].set_title(f"{key}")
10
11     table = data.loc[baseline_samples.get(key), :].copy()
12
13     # saturated or insaturated?
14     # =====
15     table = table.loc[:, [not saturated(x) and max_carbon(x, 30) for x in table.
16
17     # these need to be kept the same manually
18     if i == 0:
19         cols = table.max() > 500_000
20         cols = cols[cols].index
21     table = table[cols]
22     BASELINE_COLS = cols.copy()
23
24     bar = sns.barplot(data=table, ax=axes[i],
25                       #alpha=.5
26                       )
27     #swarm = sns.swarmplot(data=table, ax=axes[i], s=10, alpha=.85, edgecolor="
28     axes[i].tick_params(rotation=90, labelsz=10, axis="x")
29     axes[i].tick_params(labelsize=10, axis="y")
30     if i != 0:
31         axes[i].set_yticklabels([])
32     axes[i].set_ylim(0, 7_000_000)
```

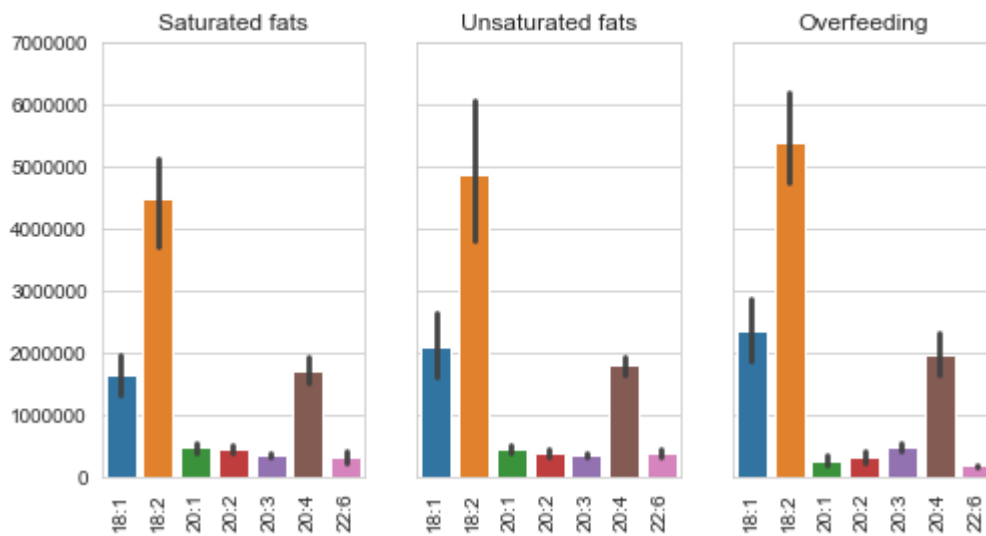
Baseline conditions, abundant, unsaturated residues



In [51]:

```
1 label = "After treatment, abundant, unsaturated residues"
2
3 fig, axes = plt.subplots(1, 3, figsize=(8,4))
4 plt.suptitle(label, size=18, y=1.05)
5 labelled_columns = 2
6
7 for i, key in enumerate(("Saturated fats", "Unsaturated fats", "Overfeeding")):
8
9     axes[i].set_title(f"{key}")
10
11     #table = data.loc[baseline_samples.get(key), :].copy()
12     table = data[data["diet"] == key].copy()
13
14     # the same as previous plot
15     table = table[BASELINE_COLS]
16
17     bar = sns.barplot(data=table, ax=axes[i],
18                       #alpha=.5
19                       )
20     #swarm = sns.swarmplot(data=table, ax=axes[i], s=10, alpha=.85, edgecolor="k")
21     axes[i].tick_params(rotation=90, labels=10, axis="x")
22     axes[i].tick_params(labels=10, axis="y")
23     if i != 0:
24         axes[i].set_yticklabels([])
25     axes[i].set_ylim(0, 7_000_000)
```

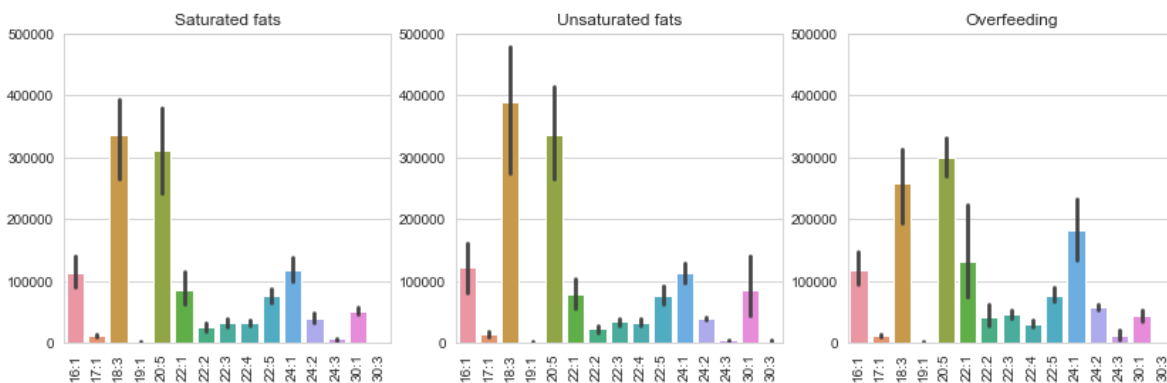
After treatment, abundant, unsaturated residues



In [52]:

```
1 label = "Baseline conditions, less abundant, unsaturated residues"
2
3 fig, axes = plt.subplots(1, 3, figsize=(14,4))
4 plt.suptitle(label, size=18, y=1.05)
5 labelled_columns = 2
6
7 for i, key in enumerate(("Saturated fats", "Unsaturated fats", "Overfeeding")):
8
9     axes[i].set_title(f"{key}")
10
11     table = data.loc[baseline_samples.get(key), :].copy()
12
13     # saturated or insaturated?
14     # =====
15     table = table.loc[:, [not saturated(x) and max_carbon(x, 30) for x in table.
16
17     # these need to be kept the same manually
18     if i == 0:
19         cols = table.max() <= 500_000
20         cols = cols[cols].index
21     table = table[cols]
22     BASELINE_COLS = cols.copy()
23
24     bar = sns.barplot(data=table, ax=axes[i],
25                       #alpha=.5
26                       )
27     #swarm = sns.swarmplot(data=table, ax=axes[i], s=10, alpha=.85, edgecolor="k")
28     axes[i].tick_params(rotation=90, labelsz=10, axis="x")
29     axes[i].tick_params(labelsize=10, axis="y")
30     axes[i].set_ylim(0, 500_000)
```

Baseline conditions, less abundant, unsaturated residues



In [53]:

```
1 label = "After treatment, less abundant, unsaturated residues"
2
3 fig, axes = plt.subplots(1, 3, figsize=(14,4))
4 plt.suptitle(label, size=18, y=1.05)
5 labelled_columns = 2
6
7 for i, key in enumerate(("Saturated fats", "Unsaturated fats", "Overfeeding")):
8
9     axes[i].set_title(f"{key}")
10
11     #table = data.loc[baseline_samples.get(key), :].copy()
12     table = data[data["diet"] == key].copy()
13
14     # saturated or insaturated?
15     # =====
16     table = table.loc[:, [not saturated(x) and max_carbon(x, 30) for x in table.
17
18     # the same as previous plot
19     table = table[BASLINE_COLS]
20
21     bar = sns.barplot(data=table, ax=axes[i],
22                       #alpha=.5
23                       )
24     #swarm = sns.swarmplot(data=table, ax=axes[i], s=10, alpha=.85, edgecolor="k")
25     axes[i].tick_params(rotation=90, labelsize=10, axis="x")
26     axes[i].tick_params(labelsize=10, axis="y")
27     axes[i].set_ylim(0, 500_000)
```

After treatment, less abundant, unsaturated residues

