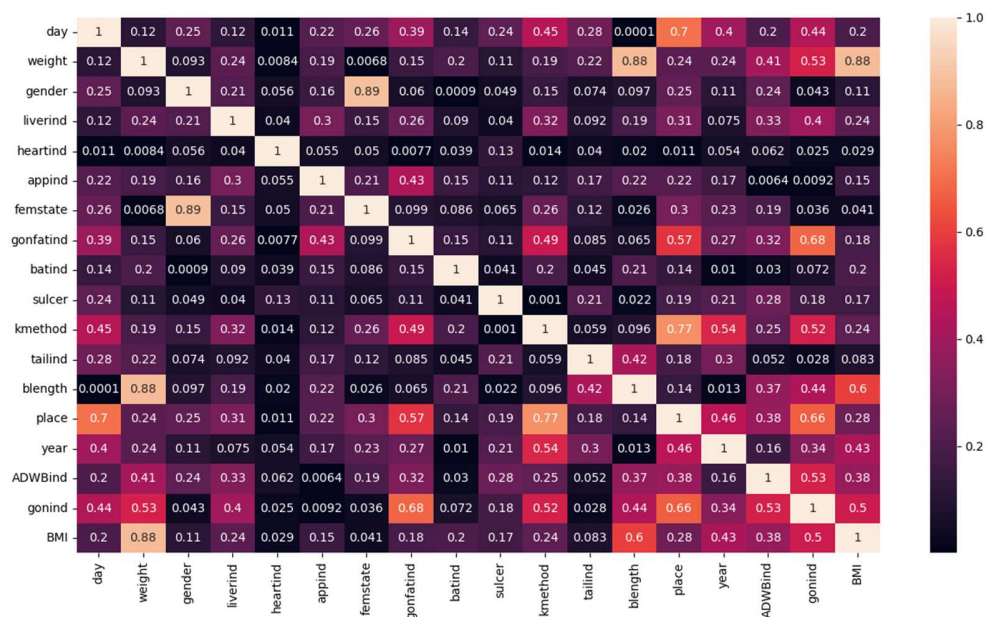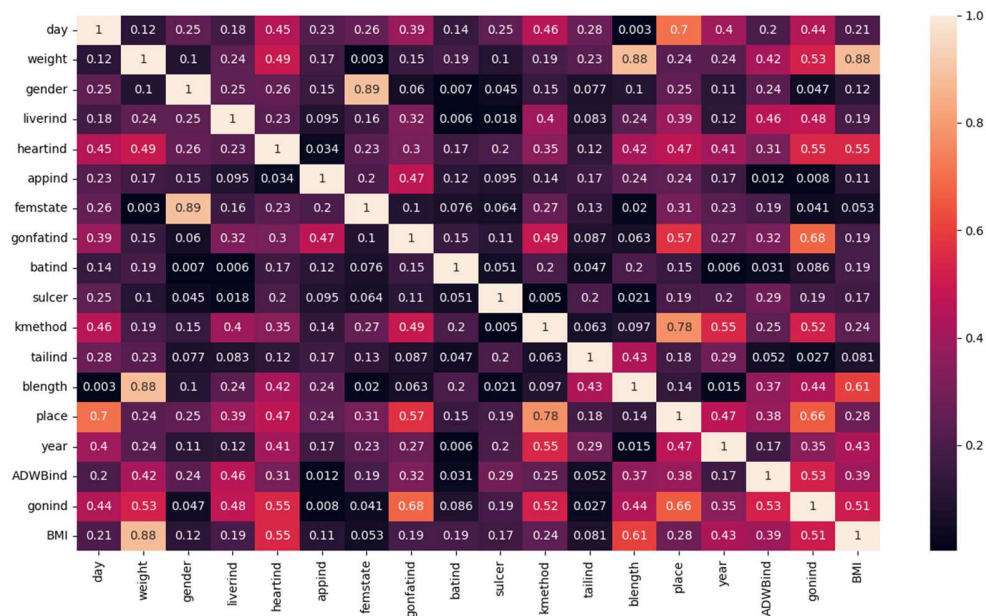Task1

a)  i) The strongest correlation involving categorical features is between **gender** and **femstate**, where the correlation is about **0.89**.

ii) The strongest correlation involving temporal features is between **day** and **place**, where the correlation is about **0.7**.

iii) The strongest correlation involving other numerical features is between **weight** and **blength**, where the correlation is about **0.88**.



b)  Also, with **liverind,** but especially with **heartind** the largest correlations increased significantly. Before the largest correlation involving heartind was 0.13, but after removing outliers the largest correlations were about 0.55. The reason for the huge difference are abnormally large or small values for the two features with the outliers in the dataset that could be due to some very special cases with the outliers or mistakes in measuring or recording the data. Ie. rat120 has a 100-times higher heartind than some of the other rats.

c)   When changing the values of **day** and **year** for freezer rats some correlations changed a lot. Removing all freezer rats made the correlations normalize back to similar numbers as before. This is mainly due to relationships between different features. Ie. the largest **place** values (7-9) all belong to freezer rats. We can say that "the larger the day the larger the place", but if we change the day of freezer rats from 400 to 0, the correlation changes. Removing all freezer rats almost normalizes the situation since there are only 34 of them and they only represent a small portion of the data.

d)   .

e)   A similar correlation as in part c can be found between **place** and **kmethod**. Ie. when kmethod is 1, 317 rats have a place value of 2. That is 85% of the rats that have 1 as kmethod. By changing each kmethod that has value 2 to 10, we can reduce the correlation between place and kmethod from 0.74 down to 0.03.

f)   The most reliable correlations are the strongest ones. Between **femstate** and **gender, weight** and **blength** and between **weight** and **BMI**, the correlations are almost 0.9 and in most cases it applies that the higher weight the higher blength and BMI, so these correlations for example show strong linear trends.

```python
import pandas as pd
import numpy as np
import math
import matplotlib.pyplot as plt
import seaborn as sns
import statistics


df = pd.read_csv('corrtestdata.csv')

#remove outliers
i = 0
indices = []
for row in df['id']:
    if row in ['rat2', 'rat53', 'rat120', 'rat434']:
        indices.append(i)
    i+=1



df.drop(indices, inplace=True)

#remove freezed
i = 0
indices = []
for row in df['day']:
    if row == 400:
        indices.append(i)
    i+=1

df.drop(indices, inplace=True)


df.loc[df['kmethod'] == 1, 'kmethod'] = 10

#remove ratID
df.drop(columns=['id'], inplace=True)

##change freezed
##df.loc[df['day'] == 400, 'day'] = 0
##df.loc[df['year'] == 15, 'year'] = -1

c = df.corr().abs().round(3)

sorted = s = c.unstack().sort_values()
```

```python
plt.figure(figsize = (15, 15))
sns.heatmap(c, annot = True)

##plt.matshow(df.corr())
plt.show()
```