

Exercise Sheet 1

In the following, we explore different ways of accessing data, including reading CSV files, querying databases, and applying preprocessing and plotting techniques to the available data. The cell below imports some libraries that are required to complete the tasks. Note that you need to install additional python libraries such as `cv2`, `torch`, `torchvision`, `matolotlib` and `sqlite3`. Some of these libraries will also be needed for the next exercise sheets.

```
In [4]: %matplotlib inline
import numpy, scipy, scipy.spatial
import matplotlib
from matplotlib import pyplot as plt
from PIL import Image
import sqlite3
import torch
import torchvision, torchvision.transforms
```

Exercise 1: Loading CSV Data (15+15 P)

In this exercise, we investigate the usage of the function `numpy.genfromtxt` to load several datasets from the UCI repository. These datasets are provided in the form of csv files in the folder `csvdata` of the homework.

(a) Using the method `numpy.genfromtxt`, load the dataset contained in the file `Wholesale customers data.csv`. In this dataset, instances (rows) are retailers, and features (columns) represent how much these retailers spend for different categories of products. Once the dataset is loaded, compute the average and median spending (over instances) for each category of products.

```
In [2]: # -----
# TODO: Replace by your code
# -----
df = numpy.genfromtxt('Wholesale customers data.csv', delimiter=',')
cols = numpy.genfromtxt('Wholesale customers data.csv', delimiter=',', names=True).dtype.names[2:]
df = df[1:,2:]
meanlist = numpy.mean(df, axis=0)
medianlist = numpy.median(df, axis=0)
print("{: >4s}".format('Mean'), "{: >21s}".format('Median'), '\n')
for i in range(8, len(cols)):
    print("{: >25s}".format(cols[i]), "{: >25s}".format(meanlist[i]), "{: >15s}".format(medianlist[i]), '\n' )
# -----
```

	Mean	Median
Fresh	12000.297727272728	8504.0
Milk	5796.265909090909	3627.0
Grocery	7951.277272727273	4755.5
Frozen	3071.931818181818	1526.0
Detergents_Paper	2881.4931818181817	816.5
Delicatessen	1524.8704545454545	965.5

(b) Using the method `numpy.genfromtxt`, load the dataset contained in the file `CortexNuclear.csv`, and use the library `matplotlib` to produce an image plot that visualizes the dataset, specifically visualize the first 30 instances that do not contain any missing value.

```
In [3]: # -----
# TODO: Replace by your code
# -----
plt.figure(figsize=(15, 15))
data = numpy.genfromtxt('CortexNuclear.csv', delimiter=',')
data_mod=data[:, 1:]
data_mod=data_mod[:, ~numpy.isnan(data_mod).all(axis=0)] #filter column with all null values
data_mod=data_mod[~numpy.isnan(data_mod).any(axis=1)] #filter rows with any null values
data_mod=data_mod[:30,:] #keeping only first 30 instances
plt.imshow(data_mod) #creating image plot
# -----
```

Out[3]: <matplotlib.image.AxesImage at 0x2783c8eb490>

Exercise 2: Querying a Database (20+20 P)

In the following, we will use the `sqlite3` package to connect to a database, and perform various join operations. The `sqlite3` package enables you to connect to a database and to perform various queries. We will consider the `chinook` database, which simulates data from a music store, relating music tracks, artists, invoices, customers, etc. Connect to the database. The database can also be downloaded from the link <https://www.sqliteutorial.net/sqlite-sample-database/>.

```
In [77]: db = sqlite3.connect('chinook.db')

The database has the following schema
```

We first consider a simple query on this database. The query is formulated in the SQL language and retrieves the duration of tracks found in that database. Once the results of the query have been obtained, we perform a very basic data analysis: computing the mean track duration.

```
In [78]: cursor = db.cursor()
query = "SELECT Milliseconds FROM tracks;"
results = numpy.array(cursor.execute(query).fetchall())
mean = results.mean()/1000.0
print(f'{"Average track duration":25s} {mean:8.3f}')

Average track duration      393.599
```

Now, we would like to perform more complex SQL queries. For a tutorial on SQL, see for example, <https://www.sqltutorial.org/>. In particular, look at Section 6 which discusses the SQL operation "INNER JOIN" and that is useful for generating outputs involving multiple tables.

(a) Apply a SQL query that extracts a table containing for all tracks their genre and their track length. Then, write code that computes for each genre (sorted alphabetically the average track length).

```
In [79]: # -----
# TODO: Replace by your code
# -----
query = """
        SELECT Name, AVG(Milliseconds) / 1000 FROM genres INNER JOIN (SELECT GenreId, Milliseconds FROM tracks)
        as tracks_genres USING(GenreId)
        GROUP BY Name
        ORDER BY Name;
    """

results = numpy.array(cursor.execute(query).fetchall())
for row in results:
    print(f'"{row[0]:25s} {float(row[1]):8.3f}"')
# -----
```

Alternative	264.059
Alternative & Punk	234.354
Blues	270.360
Bossa Nova	219.590
Classical	293.868
Comedy	1585.264
Drama	2575.284
Easy Listening	189.164
Electronica/Dance	302.986
Heavy Metal	297.453
Hip Hop/Rap	178.176
Jazz	291.755
Latin	232.859
Metal	309.749
Opera	174.813
Pop	229.034
R&B/Soul	220.067
Reggae	247.178
Rock	283.910
Rock And Roll	134.643
Sci Fi & Fantasy	2911.783
Science Fiction	2625.549
Soundtrack	244.371
TV Shows	2145.041
World	224.924

We would like to analyze the preference for music genres in different countries.

(b) Apply a SQL query that extracts for each invoice the country of the customer and the genre of the track the customer has purchased. Then, print in the form of a table the number of purchases for each country and genre.

```
In [91]: # -----
# TODO: Replace by your code
# -----
query = """
        WITH table1 AS
        (
            SELECT Name, BillingCountry, COUNT(*) as Amount FROM
            (
                SELECT GenreId, BillingCountry FROM
                (
                    (SELECT TrackId, BillingCountry FROM
                     (SELECT InvoiceId, BillingCountry FROM invoices) as t1
                     INNER JOIN invoice_items USING(InvoiceId)
                    ) as t2
                    INNER JOIN tracks USING(TrackId)
                ) AS t3
            )
            INNER JOIN Genres USING(GenreId)
            GROUP BY Name, billingcountry
        )
        SELECT Name AS Genre, BillingCountry AS Country, CASE WHEN Amount IS NULL THEN 0 ELSE Amount END AS Amount FROM
        (SELECT * FROM (SELECT DISTINCT Name FROM table1), (SELECT DISTINCT BillingCountry FROM table1)) as t1
        LEFT JOIN (SELECT Name, BillingCountry, Amount FROM table1)
        USING (Name, BillingCountry)
        ORDER BY Name, BillingCountry
    """

results = numpy.array(cursor.execute(query).fetchall())

COUNTRIES = 24
print(" " * 20, end='')
for i in range(COUNTRIES):
    print(f'"{results[i][1][:3]:3s}"', end='')
for i, row in enumerate(results):
    if not i % COUNTRIES:
        print()
        print(f'"{row[0]:20s}"', end='')
        print(f'"{int(row[2]):3.0f}"', end='')
# -----
```

	Arg	Aus	Aus	Bel	Bra	Can	Chi	Cze	Den	Fin	Fra	Ger	Hun	Ind	Ire	Ita	Net	Nor	Pol	Por	Spa	Swe	USA	Uni
Alternative	0	0	0	0	0	0	0	0	0	0	4	1	0	0	0	0	0	4	0	0	0	0	5	0
Alternative & Punk	9	0	0	14	7	36	2	9	4	2	31	13	3	11	4	11	5	2	7	5	4	6	58	9
Blues	0	1	0	0	6	4	2	1	0	0	2	14	2	3	1	4	0	1	1	0	0	15	0	
Bossa Nova	0	0	0	0	7	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	7	0	
Classical	0	0	2	0	6	5	1	0	0	0	10	0	0	2	0	2	0	5	0	0	0	8	0	
Comedy	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	
Drama	0	0	1	0	2	4	6	0	0	4	1	2	0	1	0	0	2	0	0	0	0	6	0	
Easy Listening	2	0	0	0	0	3	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	3	0	
Electronica/Dance	0	0	0	0	4	0	2	0	0	2	0	1	0	0	0	0	2	0	1	0	0	0	0	
Heavy Metal	0	3	0	0	0	0	0	0	0	0	3	0	0	0	0	0	0	0	2	0	0	4	0	
Hip Hop/Rap	0	0	0	2	3	0	0	0	2	0	0	0	0	0	0	0	0	0	0	1	0	4	3	
Jazz	2	0	2	0	0	13	0	3	0	2	11	2	0	10	3	0	0	0	1	2	2	12	91	
Latin	8	2	2	0	53	60	8	9	4	9	26	18	3	9	7	2	3	5	7	13	4	122	31	
Metal	7	8	7	1	15	40	3	6	6	2	20	25	5	8	3	1	2	0	0	11	3	7	64	
Pop	0	0	1	0	3	0	0	4	3	0	2	1	0	0	0	3	0	0	4	0	0	5	2	
R&B/Soul	0	0	4	2	3	5	0	2	0	0	0	0	5	4	0	0	0	0	0	2	0	12	2	
Reggae	0	2	0	0	6	7	0	0	0	1	1	0	0	0	0	0	0	0	0	2	0	6	5	
Rock	9	22	15	21	81	107	9	25	21	18	65	62	11	25	12	18	18	17	22	31	22	10157	37	
Rock And Roll	0	0	0	0	2	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	3	0	
Sci Fi & Fantasy	0	0	0	0	2	0	1	0	0	4	2	0	1	1	0	0	3	0	0	1	0	5	0	
Science Fiction	0	0	0	0	0	2	1	0	0	0	2	0	0	0	0	0	0	0	0	0	0	1	0	
Soundtrack	1	0	0	0	4	0	1	0	0	0	5	5	0	0	0	0	0	0	0	0	0	4	0	
TV Shows	0	0	4	0	0	1	2	8	0	0	1	3	4	1	7	0	0	0	0	1	0	114	0	
World	0	0	0	0	2	6	0	0	0	0	0	0	0	0	0	0	1	0	2	0	1	0	1	

Exercise 3: Representing Images (15+15 P)

Images are high-dimensional data. High-level concepts contained in these images are hard to extract directly from the raw pixel representation. In the following, we investigate the benefit of preprocessing the data with a neural network. Specifically, we want to see if representations built by the neural network enable to produce meaningful similarities or dissimilarities between images. We consider for this exercise a set of 10 images. The first 5 images depict geraniums, and the last 5 images depict ferrari cars.

```
In [92]: images = [Image.open(f'imagesdata/{i}.jpg') for i in range(10)]

Although the two groups are clearly distinct from a human point of view, we will show that distances computed on pixel values (i.e. treating an image as vector storing the multiple RGB pixel values) does not enable such distinction.
```

(a) Compute a matrix of pairwise Euclidean distances between images (images are resized to 100 x 100 for this task).

```
In [93]: resizedimages = [numpy.array(img.resize((100,100))).flatten() for img in images]

# -----
# TODO: Replace by your code
# -----
D = numpy.zeros((len(resizedimages), len(resizedimages)))
D.proto = scipy.spatial.distance.pdist(numpy.array(resizedimages).reshape((10,-1)), 'euclidean')

for i in range(D.shape[0]):
    for j in range(i+1, D.shape[1]):
        D[i, j] = D.proto[10 * i + j - ((i + 2) * (i + 1)) // 2]
D = D + D.T
# -----

The distance matrix can then be displayed using matplotlib.
```

In [94]: plt.imshow(D)
plt.colorbar()

Out[94]: <matplotlib.colorbar.Colorbar at 0x27a065d51f0>

We cannot see clear similarities (low distances) within each image group. This suggests that the pixel representation does not encode well the concepts we are interested in.

To address this limitation, we consider a state-of-the-art neural network called `densenet161` and available pretrained in the `torchvision` library. This neural network is composed of a feature extractor and a classification head. The feature extractor transforms image data (given as a `torch` tensor) into a tensor of activations where concepts are easier to predict.

(b) Compute the distance matrix between the images represented at the output of the densenet features extractor.

Hints: (1) The input images need to be converted to a torch tensor and the normalized using the function provided below before being fed to the neural network. (2) Note that the tensor at the output of the network can vary in shape due to the varying size of the input images. This can be addressed by applying the mean operation over the two dimensions representing the horizontal and vertical components.

```
In [96]: model = torchvision.models.densenet161(pretrained=True).features
model.eval();

def normalize(x):
    x = x - torch.Tensor([0.485, 0.456, 0.406]).reshape(1,-1,1,1)
    x = x / torch.Tensor([0.229, 0.224, 0.225]).reshape(1,-1,1,1)
    return x

# -----
# TODO: Replace by your code
# -----
transform = torchvision.transforms.Compose([torchvision.transforms.ToTensor()])
image_tensors = [normalize(transform(img)) for img in images]
output_images = [model(tensor) for tensor in image_tensors]
nn_images = numpy.zeros((10, len(torch.mean(output_images[0], 3).flatten().cpu().detach().numpy()))))
for i in range(10):
    nn_images[i, :] = torch.mean(output_images[i], 3).flatten().cpu().detach().numpy()

D = numpy.zeros((10,10))
D.proto = scipy.spatial.distance.pdist(numpy.array(nn_images).reshape((10,-1)), 'euclidean')

for i in range(D.shape[0]):
    for j in range(i+1, D.shape[1]):
        D[i, j] = D.proto[10 * i + j - ((i + 2) * (i + 1)) // 2]
D = D + D.T
# -----

plt.imshow(D)
plt.colorbar()
```

Out[96]: <matplotlib.colorbar.Colorbar at 0x27b0a78f160>

We observe that distances now form a block structure, where the first 5 images are clearly mutually similar, and similarly for the last 5 images.