

ICT for Health Laboratory # 4 Detect activities

Monica Visintin

Politecnico di Torino



December 9th 2021

Table of Contents

1 Description of laboratory # 4

2 View the data

3 Preliminary analysis

4 What do you have to do?

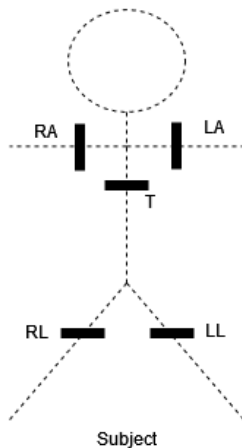
Description [1]

In fitness, people would like to automatically detect their current activity, in order to have a more precise evaluation of the calories, etc.

Task of this lab is to detect the activity (one among 19) starting from the samples taken at frequency 25 Hz from 3-axes accelerometers, gyroscopes, magnetometers placed on torso, right arm, left arm, right leg, left leg: 5 positions, 3 sensors per position, each measuring 3 values on the x-y-z coordinates, total of 45 values for each sample (<https://archive.ics.uci.edu/ml/datasets/Daily+and+Sports+Activities>).

Data were taken from 8 subjects; each subject practised the activity for 5 minutes and slices/files of 5 seconds were created (each slice/file contains 125 rows and 45 columns).

Description [2]



— 3-axis
accelerometer
gyroscope
magnetometer
(9 sensors)

Description [3]

We want a simple algorithm (not neural networks, taught in the parallel class), valid for a specific subject. The idea is to use the clustering algorithm **K-Means**:

- Read the first M slices for the 19 activities of the chosen subject.
- Apply K-Means with 19 clusters and find the 19 centroids.
- For each row of the subsequent slices detect the activity as the one with centroid at **minimum distance** from the row.

This lab is used just as a **feasibility study**. In the real life, an Android/iOS application should be written, that runs on the subject smartphone and gets the data from sensors and makes decisions. The application should have a first training phase in which the subject stores the data, labels them with the activity and applies K-Means. Note that we use a **clustering algorithm** to solve a **classification problem**.

Description [4]

Less than 15 sensors/45 features can be used (the fewer, the better).

Data can/should be **processed** before application of K-Means (for example, through smoothing, enhancing peaks, reshaping, etc) .

You cannot use machine learning algorithms other than K-Means, because simplicity is essential (less energy to perform the computations, batteries last longer, less required CPU, etc).

You can make a decision at each timestamp (25 decisions per second) or every few timestamps (for example one decision per second is enough).

The maximum number of slices in the training phase is 30 (out of a total of 60 slices); the fewer slices in the training phase, the better.

Each student will work with one of the 8 subjects:

```
1 ID = 313131 #your ID number
2 s=ID%8+1 # subject index from 1 to 8
```

Description [5]

You are responsible of your choices. **Report due by Wednesday January 17th 2023, 11.59 PM.**

Table of Contents

1 Description of laboratory # 4

2 **View the data**

3 Preliminary analysis

4 What do you have to do?

Read the data [1]

- In DropBox you find the zipped file with the time series. You can download the data from <https://archive.ics.uci.edu/ml/datasets/Daily+and+Sports+Activities>.
- Copy the file on your PC, extract the folders and files.
- **NOTE: do NOT include the data when you upload the report and the Python code.**
- In DropBox you also find a Python script that allows you to generate a Pandas dataframe that includes the selected slices of the selected activities of the selected subject.

Read the data [2]

```
1 def generateDF( filedir , colnames , patients , activities , slices ):  
2     # get the data from files for the selected patients  
3     # and selected activities  
4     # concatenate all the slices  
5     # generate a pandas dataframe with an added column: activity  
6     x=pd.DataFrame()  
7     for pat in patients:  
8         for a in activities:  
9             subdir='a'+f"{a:02d}"+' /p'+str(pat)+' / '  
10            for s in slices:  
11                filename=filedir+subdir+'s'+f"{s:02d}"+' .txt '  
12                x1=pd.read_csv( filename , names=colnames )  
13                x1['activity']=a*np.ones((x1.shape[0],), dtype=int)  
14                x=pd.concat([x,x1], axis=0, join='outer',  
15                            ignore_index=True,  
16                            keys=None, levels=None, names=None,  
17                            verify_integrity=False,  
18                            sort=False, copy=True)  
19     return x
```

Read the data [3]

● Activities

```
1 actNames=[
2     'sitting', # 1
3     'standing', # 2
4     'lying on back', # 3
5     'lying on right side', # 4
6     'ascending stairs', # 5
7     'descending stairs', # 6
8     'standing in an elevator still', # 7
9     'moving around in an elevator', # 8
10    'walking in a parking lot', # 9
11    'walking on a treadmill with a speed of 4 km/h in flat', # 10
12    'walking on a treadmill with a speed of 4 km/h in 15 deg inclined position', # 11
13    'running on a treadmill with a speed of 8 km/h', # 12
14    'exercising on a stepper', # 13
15    'exercising on a cross trainer', # 14
16    'cycling on an exercise bike in horizontal positions', # 15
17    'cycling on an exercise bike in vertical positions', # 16
18    'rowing', # 17
19    'jumping', # 18
20    'playing basketball' # 19
21 ]
```

Read the data [4]

- Short names for activities:

```
1 actNamesShort=[
2     'sitting', # 1
3     'standing', # 2
4     'lying.ba', # 3
5     'lying.ri', # 4
6     'asc.sta', # 5
7     'desc.sta', # 6
8     'stand.elev', # 7
9     'mov.elev', # 8
10    'walk.park', # 9
11    'walk.4.fl', # 10
12    'walk.4.15', # 11
13    'run.8', # 12
14    'exer.step', # 13
15    'exer.train', # 14
16    'cycl.hor', # 15
17    'cycl.ver', # 16
18    'rowing', # 17
19    'jumping', # 18
20    'play.bb', # 19
21    ]
```

Read the data [5]

● Sensors

```
1 sensNames=[
2     'T_xacc', 'T_yacc', 'T_zacc',
3     'T_xgyro', 'T_ygyro', 'T_zgyro',
4     'T_xmag', 'T_ymag', 'T_zmag',
5     'RA_xacc', 'RA_yacc', 'RA_zacc',
6     'RA_xgyro', 'RA_ygyro', 'RA_zgyro',
7     'RA_xmag', 'RA_ymag', 'RA_zmag',
8     'LA_xacc', 'LA_yacc', 'LA_zacc',
9     'LA_xgyro', 'LA_ygyro', 'LA_zgyro',
10    'LA_xmag', 'LA_ymag', 'LA_zmag',
11    'RL_xacc', 'RL_yacc', 'RL_zacc',
12    'RL_xgyro', 'RL_ygyro', 'RL_zgyro',
13    'RL_xmag', 'RL_ymag', 'RL_zmag',
14    'LL_xacc', 'LL_yacc', 'LL_zacc',
15    'LL_xgyro', 'LL_ygyro', 'LL_zgyro',
16    'LL_xmag', 'LL_ymag', 'LL_zmag']
```

Read the data [6]

- Folder with data:

```
1 file_dir = '../data/'
```

In the given version, a folder e.g. lab3 includes folder data and folder surname_s313131_report_2_code; folder surname_s313131_report_2_code includes file data_visualization_forstudents.py

Read the data [7]

- Settings:

```
1 ID=313131
2 s=ID%8+1
3 patients=[s] # list of selected patients
4 activities=list(range(1,6)) #list of indexes of activities to plot
5 Num_activities=len(activities)
6 NAc=19 # total number of activities
7 actNamesSub=[actNamesShort[i-1] for i in activities] # short names of the s
8 sensors=list(range(9)) # list of sensors
9 sensNamesSub=[sensNames[i] for i in sensors] # names of selected sensors
10 Nslices=12 # number of slices to plot
11 #Ntot=60 #total number of slices
12 slices=list(range(1,Nslices+1)) # first Nslices to plot
13 fs=25 # Hz, sampling frequency
14 samplesPerSlice=fs*5 # samples in each slice
```

In the example only the first 9 sensors and the first 5 activities are considered. In plots only the first 12 slices (60 seconds) are shown. **It is highly recommended that you see the time series associated with all the activities and possibly all the sensors.**

1 Description of laboratory # 4

2 View the data

3 Preliminary analysis

4 What do you have to do?

“Cluster” analysis [1]

- In this lab we know what we would like K-Means to give us: each cluster should include exactly all the rows/samples of one and only one activity.
- It is therefore useful to start analyzing the data as they are, assuming that K-Means will then be able to correctly find the clusters.
- The cluster centroids can be obtained as the mean of the data for each activity. It is also interesting to find the standard deviations of the measures of each sensor for each activity, because this tells us how much the the points move around the mean (for each sensor/feature).

“Cluster” analysis [2]

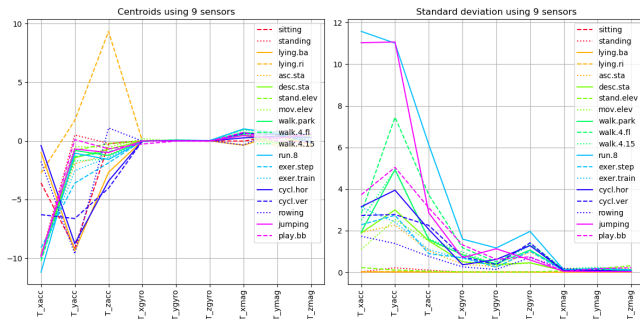
All the activities are here considered, but only 9 features are used:

```
1  ##### plot centroids and stand. dev. of sensor values
2  print('Number of used sensors: ', len(sensors))
3  centroids=np.zeros((NAc, len(sensors)))# centroids for all the activities
4  stdpoints=np.zeros((NAc, len(sensors)))# variance in cluster for each sensor
5  plt.figure(figsize=(12,6))
6  for i in range(1,NAc+1):
7      activities=[i]
8      x=generateDF( filedir ,sensNamesSub ,patients ,activities ,slices)
9      x=x.drop(columns=[' activity '])
10     centroids[i-1,:]=x.mean().values
11     plt.subplot(1,2,1)
12     lines = plt.plot(centroids[i-1,:],label=actNamesShort[i-1])
13     lines[0].set_color(cm(i//3*3/NAc))
14     lines[0].set_linestyle(line_styles[i%3])
15     stdpoints[i-1]=np.sqrt(x.var().values)
16     plt.subplot(1,2,2)
17     lines = plt.plot(stdpoints[i-1,:],label=actNamesShort[i-1])
18     lines[0].set_color(cm(i//3*3/NAc))
19     lines[0].set_linestyle(line_styles[i%3])
20 plt.subplot(1,2,1)
21 plt.legend(loc='upper right')
22 plt.grid()
```

“Cluster” analysis [3]

```
23 plt.title('Centroids using '+str(len(sensors))+ ' sensors')
24 plt.xticks(np.arange(x.shape[1]), list(x.columns), rotation=90)
25 plt.subplot(1,2,2)
26 plt.legend(loc='upper right')
27 plt.grid()
28 plt.title('Standard deviation using '+str(len(sensors))+ ' sensors')
29 plt.xticks(np.arange(x.shape[1]), list(x.columns), rotation=90)
30 plt.tight_layout()
```

“Cluster” analysis [4]



It seems that jumping and running at 8 km/h imply large variations of the values measured by the torso accelerometer, whereas the torso gyroscope mean and standard deviation values seem not to depend that much on the activity.

“Cluster” analysis [5]

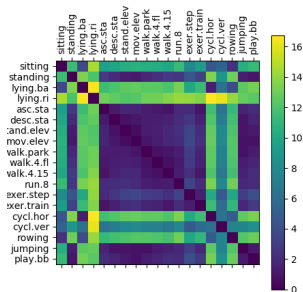
- Let us now find the distance among the centroids. The centroids should be well separated, so that the subsequent classification can be more easily and correctly accomplished.

```

1  ### between centroids distance
2  d=np.zeros((NAc,NAc))
3  for i in range(NAc):
4      for j in range(NAc):
5          d[i,j]=np.linalg.norm(centroids[i]-centroids[j])
6
7  plt.matshow(d)
8  plt.colorbar()
9  plt.xticks(np.arange(NAc),actNamesShort,rotation=90)
10 plt.yticks(np.arange(NAc),actNamesShort)

```

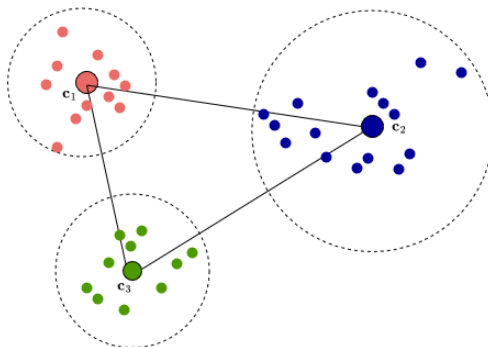
“Cluster” analysis [6]



It seems that it will be very difficult to distinguish the activities from ascending the stairs to exercising on a cross trainer, at least using the data as they are and just 9 features, because the distances among these centroids are very small.

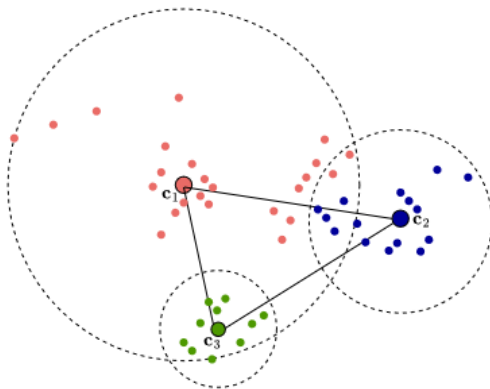
“Cluster” analysis [7]

Desired case: clusters are far apart and points are close to their centroids. K-Means can work correctly.



“Cluster” analysis [8]

Undesired case: points are very spread and distant from centroids.
K-Means cannot work correctly (it finds other centroids and clusters).



“Cluster” analysis [9]

- For each activity find the minimum distance between the corresponding centroid and the other centroids:

```
1 dd=d+np.eye(NAc)*1e6# remove zeros on the diagonal (distance of centroid from  
2 dmin=dd.min(axis=0)# find the minimum distance for each centroid
```

- Find the average distance of the points of a cluster from their centroid.
The average square distance is

$$d^2 = \frac{1}{N} \sum_{n=1}^N \sum_{f=1}^F (\mathbf{x}_{nf} - \mathbf{c}_f)^2$$

where N is the number of points in the cluster, F is the number of features/columns, \mathbf{c} is the centroid, \mathbf{x}_n is the n -th point in the cluster, \mathbf{x}_{nf} is the value of the f -th feature of the n -th point. Note that

$$d^2 = \sum_{f=1}^F \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_{nf} - \mathbf{c}_f)^2 = \sum_{f=1}^F \text{var}_f$$

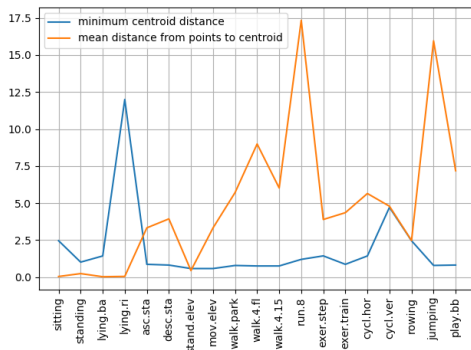
“Cluster” analysis [10]

where var_f is the variance of feature f in the cluster. Therefore, the average distance can be evaluated as

$$d = \sqrt{\sum_{f=1}^F \text{var}_f}$$

```
1 dpoints=np.sqrt(np.sum(stdpoints**2,axis=1))
```

“Cluster” analysis [11]



If the mean distance of points from the centroid is higher than the minimum distance between the centroid and the other centroids, then we are in the undesired case and the corresponding activity cannot be detected. Only the first 4 activities can be easily detected.

- 1 Description of laboratory # 4
- 2 View the data
- 3 Preliminary analysis
- 4 **What do you have to do?**

To do [1]

- 1 For your specific subject (that depends on your ID number), have a look at all the plots, all the activities, possibly all the features/sensors/axis.
- 2 Work on the time series/data, through data processing, to improve the plot in slide 26 (mean distance of points in the cluster should be less than the minimum distance between the cluster centroid and all the other centroids).
- 3 Select the number of training slices, select **all the activities** and generate the training and test datasets. If, for example, you want to use 10 slices for training, then your training dataset will have $10 \times 125 \times 19$ rows, the test dataset will have $50 \times 125 \times 19$ rows (at least with the current initial settings, with no data processing). It is not necessary to shuffle the data, since K-Means does not take advantages of data order.
- 4 Run Scikit-Learn K-Means algorithm (fit the training dataset) and run it on the test dataset (predict). Choose appropriate parameters in the Scikit-Learn K-Means call.

To do [2]

- 5 **Note:** the **index of clusters at the output of K-Means is random**. Not necessarily activity 1 corresponds to cluster 1, activity 1 might correspond to cluster 10. Find the **mapping** between K-Means cluster indexes and activities, so that you can correctly remap cluster indexes to activities.
- 6 For both training and test datasets, find the confusion matrix that gives

$$P(i|j) = P(\text{detected activity is } i | \text{the true activity is } j)$$

in row j column i , for $i = 1, \dots, 19$ and $j = 1, \dots, 19$. Index i corresponds to activity i , index j corresponds to activity j .

To do [3]

- For both training and test datasets, find the accuracy:

$$\text{Accuracy} = \sum_{i=1}^{19} P(i|i)P(i)$$

where $P(i)$ is the probability of activity i , $P(i) = 1/19$. These values must be written in the report.

- Find the accuracy of each activity and write the corresponding table in the report. Find if two activities are easily confused, etc.
- With an appropriate data processing, you should get an accuracy around 0.85-0.95 for both training and test data, these values change with the subject. If you cannot, try with other time series processing techniques. If you cannot find a solution turn in your best solution.
- Write the report, **6 pages** (same rules as for report 1), in which you exactly explain what you did, so that another researcher can repeat your experiment and obtain the same results.

To do [4]

- 11 The report must include at least the plot of the confusion matrix, the accuracies for each of the 19 activities, the overall accuracy for training and test datasets.
- 12 **NOTE: do NOT include the data when you upload the report and the Python code.**
- 13 It is forbidden to **copy** figures or sentences from these slides, from the lecture slides, from Wikipedia, from the web, etc (this is **plagiarism**). You must write your own sentences with your own words. If you need a figure, you produce it.
- 14 You **cannot discuss** the method you decided to adopt with the other students, each student must work on his/her own. If the methods (and/or the report) of two students are too similar, the grade will be divided by 2 (if you find a solution, don't tell it to other students otherwise you'll get half of the correct grade, do not trust your colleagues).

To do [5]

Start being men and women responsible of their actions, not cheating children. You do not learn if you cheat, and you are paying for learning (both time and money). You will not be able to cheat in working activities: your colleagues will have their own job to do, they will not help you.