16/6/20                   Assignment 08                    Yuval S. Katz
                                                            204025258

1) a- ex8_sec1.png + ass_1.py.

   b- PC1 dimension spans 2 main clusters around [0-30] and [-40,-10].
      with 2 centers. +dense prob. (0,0.5)
      PC2 dimension spans data around [0-15] with no clear division.
   c- PC1 explains 58.3%. } of the variance, ═⟩ of the variance in the data
      PC2 explains 37.2%. } together 95.5% ═⟩ can be explained using PC1+2
                                                                    PC1+2
      (out of 20)        1-95.5% remains unexplained by this representation
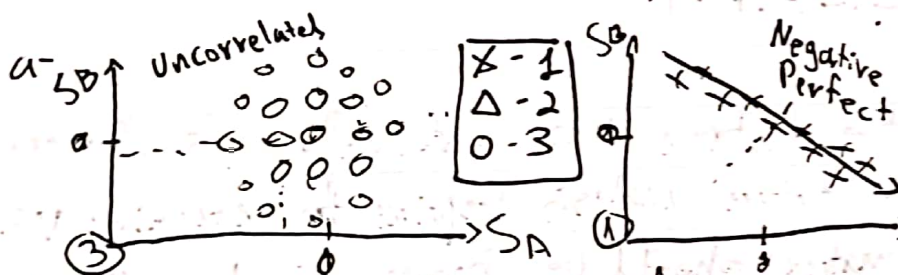   d- Vector 13 has max_var, with entropy of 2.84 bits on avg to pass information
      PC1 explains has high variance, with    "    "   2.83    "   .
      I used 10 bins for the dist. calculations. Since PC1 gets most of the
      attribute of the data orthogonaly to the features, it uses a similar dist.
                              (dimension)
      to represent the data, as do vector /feature 13. thus, the same
      (or at least very close) entropy // (regardless of bin size, just to represent
      the data to a suficient level).

2) a-

   $cov_1(S_A, S_B) = -1$
   $cov_2(S_A, S_B) = 0.5$
   $cov_3(S_A, S_B) = 0$

   b- ass_2.py
      Simulation of values with 1-3 cov:

      ① redundent! the same parameter
      ② λ1=15 λ2=0.1 V1=[0.7...] [0.7...]
      ③ redundent! V1 M=1

   c- The first (1) is redundent because $S_B$ dosent add information beyond
      $S_A$ information. (cov=-1). it sufficient to look at either $S_A$/$S_B$ alone!
      The last (cov=0,0) display the highest redundancy becase $S_A$
      and $S_B$ are unrelated (=uncorrelated) and hence, adds a lot of inform.

3) a- see file ex8_ass3.py + .png.

b- on the fig ___

c- on the fig -)    red    black    purple    blue    green    } After 60
         $(-7,-3)$ $(-5,5)$ $(0,-7)$ $(5,-7)$ $(6,3)$ } iterations
centroids

4) EM 6m.

a- see file ass_4.py + .png

b- $\mu_i, \sigma_i \Longrightarrow$   $\mu_1, \sigma_1$ , $\mu_2, \sigma_2$   $\mu_3 \sigma_3$   $\mu_4 \sigma_4$
    $k = i$    $(4.35, 2.8)$ $(1.05, 0.87)$ $(4.075, 0.7)$ $(3.9, 4.83)$

c- on fig

d- on fig ass_4_3.png

e) That depends on the distance matrix. Visually, k=3 clusters fit
the data well, but future analysis requires to determine vs. k=4,
and a distance matrix should be chosen. (Eclidian /Manhatan....)
for example, we can make a .DBSCAN and look for outliers...