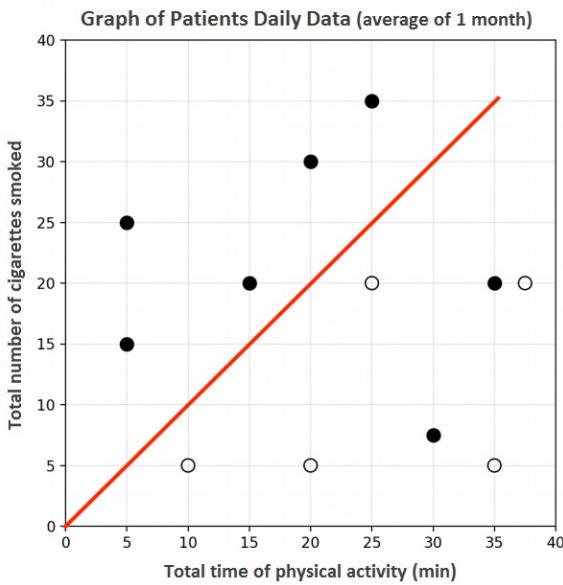


question 1)

a)



b) first most nearest point to (5,21) is (5,20) which has lung cancer. then this point is classified as with lung cancer.

c) first of all we transform the data to new space:

x	y	label	new space
5	15	1	$0.7 * (20)$
5	25	1	$0.7 * (30)$
15	20	1	$0.7 * 35$
20	30	1	$0.7 * 50$
25	35	1	$0.7 * 60$
30	7.5	1	$0.7 * 37.5$
35	20	1	$0.7 * 55$
10	5	0	$0.7 * 15$
20	5	0	$0.7 * 25$
25	20	0	$0.7 * 45$

x	y	label	new space
35	5	0	$0.7 * 40$
37.5	20	0	$0.7 * 57.5$
5	21	0	$0.7 * 26$

label of the most nearest point to the given point is healthy. Then this point is classified as healthy.

d) No. because in PCA data is transformed to the new space without considering label of the points and projects them to the direction of maximum variance. That means it doesn't project data to a new space where different classes are separated well from each other.

e, f)

	gene 1	gene 2	gene 3
mean	1.3	6	7.6
var	9.21	0.8	0.840

var is calculated as follows:

$$\text{var}(x_1) = 1/10 * (1 + 36 + 1 + 4 + 0 + 1 + 1 + 4 + 25 + 36) - 1.3 * 1.3$$

$$\text{var}(x_2) = 1/10 * (36 + 36 + 25 + 36 + 64 + 25 + 36 + 25 + 36 + 49) - 36$$

$$\text{var}(x_3) = 1/10 * (49 + 81 + 49 + 49 + 64 + 81 + 49 + 64 + 64 + 36) - 7.6 * 7.6$$

g and h)

g) \bar{X} is:

$$\begin{bmatrix} -r_1r & \varepsilon_1v & -r_1r & \varepsilon_1v & -1, r & -r_1r & -r_1r & -r_1r & r_1v & \varepsilon_1v \\ 0 & 0 & -1 & 0 & 1 & -1 & 0 & -1 & 0 & 1 \\ -, 9 & 1, 2 & -, 9 & -, 9 & -, 8 & 1, 8 & -, 9 & -, 8 & 1, 2 & -, 9 \end{bmatrix}$$

h) $C = \begin{bmatrix} 9, 21 & 1 & -, 1 \\ 1 & -, 1 & -, 1 \\ -, 1 & -, 1 & 1, 2 \end{bmatrix}$

We must calculate $\bar{X}\bar{X}^T \cdot \frac{1}{n}$
for filling this matrix

$$C_{11} = \frac{(-r_1r \times -1 + r_1r \times -1, r - 1 \times r_1v - 1 \times -r_1r + 1 \times \varepsilon_1v)}{18} = 1$$

$$C_{12} = \frac{1}{18} \times (-r_1r \times -, 9 + \varepsilon_1v \times 1, 2 - r_1r \times -, 9 + \varepsilon_1v \times -, 9) = -0, 1$$

$$= -0, 1$$

$$C_{23} = \frac{1}{18} (0, 4 \times -, 1 - 1, 2 \times -, 8 - 1, 2 \times -, 9) = -0, 1$$

~~2 4 2 0 0 0 5 0 0~~

~~T 1 1 1 1 1 1 1 1~~

~~1 1 1 1 1 1 1 1 1~~

~~1 1 1 1 1 1 1 1 1~~

~~1 1 1 1 1 1 1 1 1~~

C_{11}, C_{22}, C_{33} are calculated in the previous part

I and J)

$$\text{I) } \nu c = \nu \lambda \rightarrow \nu c = \nu \lambda I \rightarrow \nu (c - \lambda I) = 0$$

$$c - \lambda I = \begin{bmatrix} 9.121 - \lambda & 1 & -0.101 \\ 1 & -0.1 - \lambda & -0.1^2 \\ -0.101 & -0.1^2 & 0.12 - \lambda \end{bmatrix}$$

$$|c - \lambda I| = (9.121 - \lambda)[(-0.1 - \lambda)(0.12 - \lambda) + 0.02]$$

$$+ 1 \times \left(\frac{14}{1000} - \frac{0.02}{100} \times \lambda \right) - \frac{1}{100} \left[-0.1 - \frac{1}{100} (0.1 - \lambda) \right] = 0$$

after solving the above eq:

$$\lambda_1 = 9.121 \quad \lambda_2 = 0.1940 \quad \lambda_3 = 0.1004$$

$$\text{J) } \frac{9.121}{9.121 + 0.1940 + 0.1004} = \frac{9.121}{10.4154} = 0.871\%$$

K)

$$K \begin{pmatrix} 9.1\omega & 1 & -1.1\omega \\ 1 & 2\omega & -0.1\omega \\ -1.1\omega & -0.1\omega & 1.1\omega \end{pmatrix} \begin{pmatrix} v_1 \\ v_2 \\ v_3 \end{pmatrix} = \lambda \begin{pmatrix} v_1 \\ v_2 \\ v_3 \end{pmatrix}$$

$$9.1\omega v_1 + v_2 - 1.1\omega v_3 = 9.1\omega^2 v_1 \quad (1)$$

$$v_1 + 2\omega v_2 - 0.1\omega v_3 = 2\omega^2 v_1 \quad (2)$$

$$-1.1\omega v_1 - 0.1\omega v_2 + 1.1\omega v_3 = 1.1\omega^2 v_3 \quad (3)$$

by solving above eq's we obtain v as follows:

$$(1) \quad v = \begin{bmatrix} 0.99\omega \\ 0.114\omega \\ -0.100\omega \end{bmatrix}, \quad (2) \quad v = \begin{bmatrix} 0.10\omega \\ -0.100\omega \\ 0.1\omega \end{bmatrix}, \quad (3) \quad v = \begin{bmatrix} -0.9 \\ -0.11\omega \\ -0.100\omega \end{bmatrix}$$

$$\Rightarrow v = \begin{bmatrix} 0.99\omega & 0.10\omega & -0.9 \\ 0.114\omega & -0.100\omega & -0.11\omega \\ -0.100\omega & 0.1\omega & -0.100\omega \end{bmatrix}$$

L)

$$L) Y = V^t \bar{X} \quad V^t = 3 \times 3, \bar{X} = 3 \times 10 \rightarrow V^t \bar{X} = 3 \times 10$$

$\begin{bmatrix} 1.12V & -0.49 & 0.29 & -0.16 & 1.0V & 1.22 & 0.21 & 0.15 \\ 0.4V & -1.01 & 0.119 & 0.18 & 0.1V & -1.02 & 0.4V & -0.19 \\ -0.15 & 0.15 & -0.902 & 0.15 & 1.19 & 0.14 & -0.114 & -0.171 \end{bmatrix}$

$\begin{bmatrix} -0.1418 & 1.07 \\ -0.110 & -0.189 \end{bmatrix}$

question 2)

A and B and C)

q) transformation matrix V is: $S_w (M_1 - M_2)$

$$M_1 - M_2 = \begin{bmatrix} -r_1v \\ -\epsilon_1\omega \end{bmatrix}$$
$$S_w^{-1} \rightarrow S_w = S_1 \cap S_2$$
$$S_w = \begin{bmatrix} r_1v & r_1v\omega \\ r_1v\omega & \epsilon_1\omega \end{bmatrix}, S_w^{-1} = \frac{1}{|\text{Det}|} \begin{bmatrix} \epsilon_1\omega & -r_1v\omega \\ -r_1v\omega & r_1v^2 \end{bmatrix}$$
$$|S_w| = r_1v \times \epsilon_1\omega - r_1v\omega \times r_1v\omega = 1 \cdot r_1\omega - v_1 r_1 \omega \epsilon_1 \omega = r_1 \omega \epsilon_1 \omega$$
$$= r_1 \omega \epsilon_1 \omega$$
$$S_w^{-1} = \frac{1}{r_1 \omega \epsilon_1 \omega} \times \begin{bmatrix} \epsilon_1\omega & -r_1v\omega \\ -r_1v\omega & r_1v^2 \end{bmatrix} = \begin{bmatrix} 1/\omega & -1/\omega \\ -1/\omega & 1/v^2 \end{bmatrix}$$
$$S_w(M_1 - M_2) = \begin{bmatrix} -1/\omega \\ 1/v^2 \end{bmatrix}$$

discriminant func is: $y = -1/\omega r_1 + 1/v^2 \omega$

b) $y = -1/\omega r_1 + 1/v^2 \omega \times \epsilon_1 \omega = -1/v^2 \omega$

center of first class: $\mu_1 = -1,95 \times 1,1 + -1,80 \times 1,19 = -1,90,1$

center of second class: $\mu_2 = -1,95 \times 0,0 + -1,80 \times 1,01 = -1,91,9$

$$|-1,91,9 + 1,90,1| < |-1,91,9 + 1,91,9|$$

\uparrow_{new} \uparrow_{new}

μ_1 μ_2

→ this point belongs to the first class.

c) n independent variables are normally distributed

the same across levels of predictors.

3) independence: participants are randomly sampled.

assumptions 1 and 3 are not clearly declared in the

question, but assumption 2 is valid according to

the given scatters.

D)

d) by this new scatter matrix, the assumption of homogeneity will fail. Then it would be reasonable to use quadratic discriminant analysis instead of PCA.

e)

F and G)

f)

g) $M_1 = [w, w_4]$ $M_2 = [v, v_1, F]$

$X_1 M_1 = \begin{bmatrix} 0 & -1 & 1 & -1 & 1 \\ 1 & 1 & -1 & 1 & -1 \end{bmatrix} \rightarrow S_1 = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$

$X_2 M_2 = \begin{bmatrix} 1 & -1 & 0 & 1 & -1 \\ 1 & 1 & -1 & 1 & -1 \end{bmatrix} \rightarrow S_2 = \begin{bmatrix} 1 & 0 \\ -1 & 1 \end{bmatrix}$

$S_w = S_1 + S_2 = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}, S_w^{-1} = \frac{1}{(1 \cdot 1 - 1 \cdot -1)} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}$

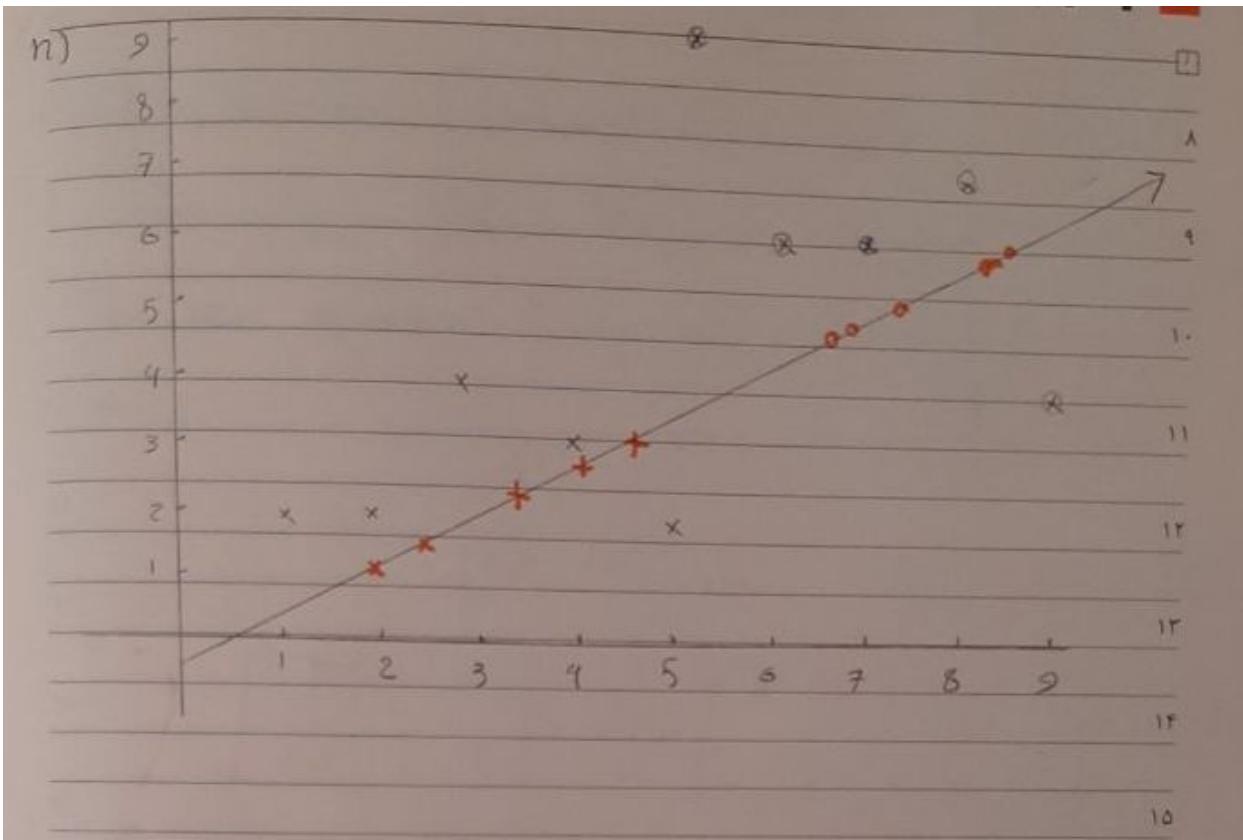
$S_w = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix}, V = S_w^{-1} (M_1 - M_2) = S_w^{-1} [-1, -1] =$

$V = \begin{bmatrix} -1 & -1 \\ -1 & -1 \end{bmatrix}$

$V = \begin{bmatrix} -1 & -1 \\ -1 & -1 \end{bmatrix} \rightarrow Y = -1 w_1 M_1 - 1 v_1 M_2$

C D E F G
T 1 T 1
A V F D F T
B 10 14 13 12 11 10
T 12 11 10 13 12 14
T 13 12 11 10 14

H and I)



l) As we see in the above figure two classes are well separated in new subspace.

new coordinate of points in above figure

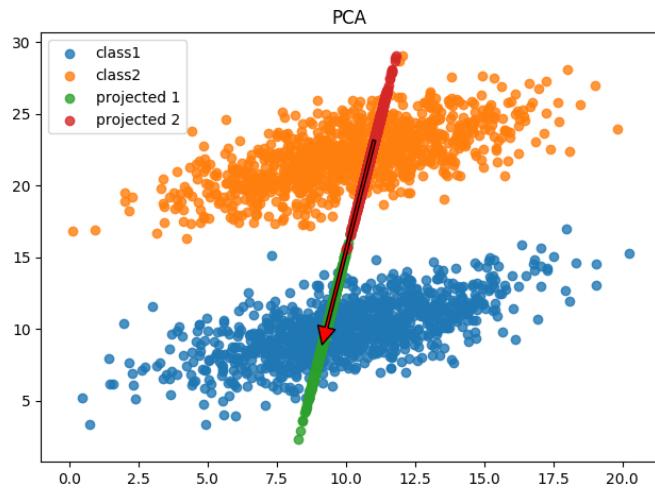
$$\begin{bmatrix} v_{1,1} & v_{1,2} & v_{1,3} & v_{1,4} \\ v_{1,5} & v_{1,6} & v_{1,7} & v_{1,8} \end{bmatrix} \quad \begin{bmatrix} v_{2,1} & v_{2,2} & v_{2,3} & v_{2,4} \\ v_{2,5} & v_{2,6} & v_{2,7} & v_{2,8} \end{bmatrix}$$

class one

class two

question 3)

a, b, c)



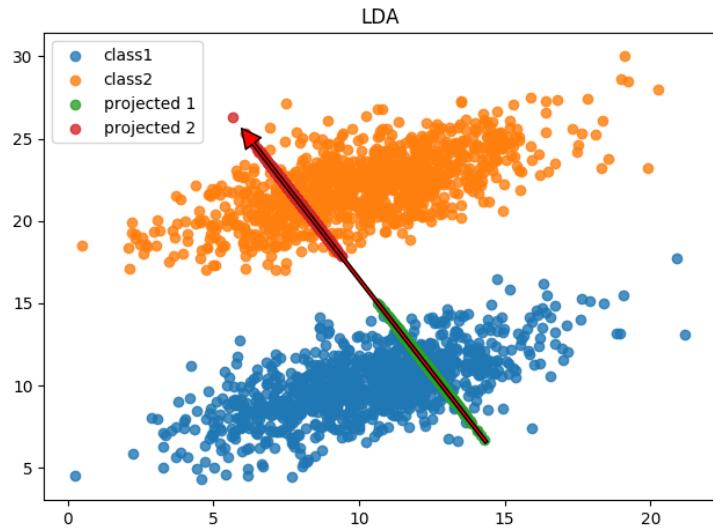
d) We expect that PCA projects the data to the line in direction of maximum variance, and as we see in the above plot, PCA has founded a projection line which is in the direction of maximum variance.

e) reconstruction error is obtained by following formula (MSE):

```
np.sqrt(np.sum([[x**2+y**2] for x,y in res]))
```

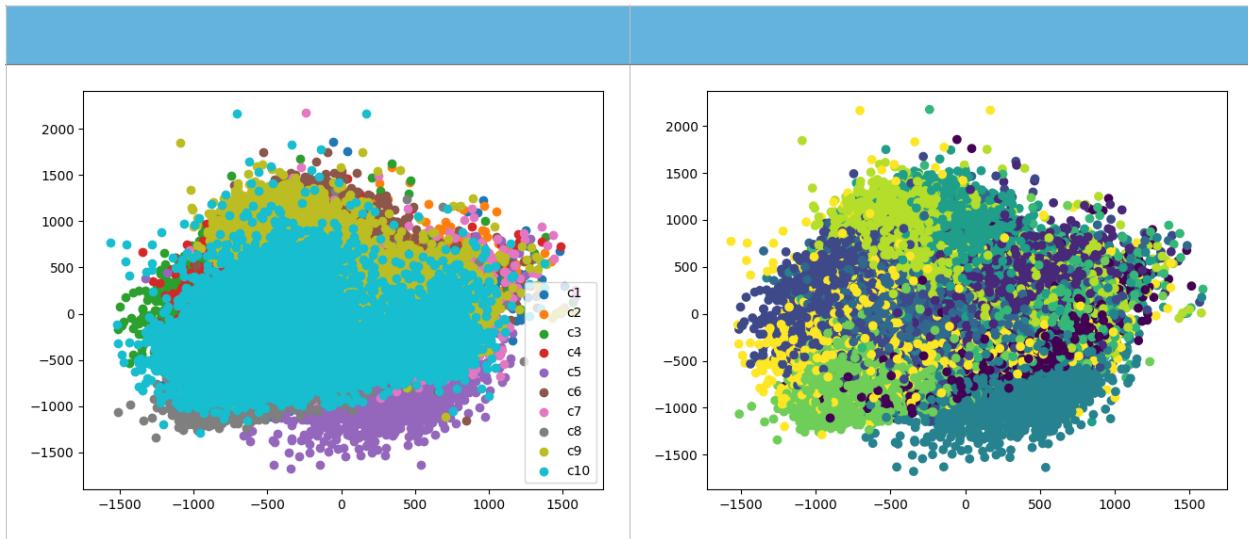
and is equal to: 17.2796

f, g)

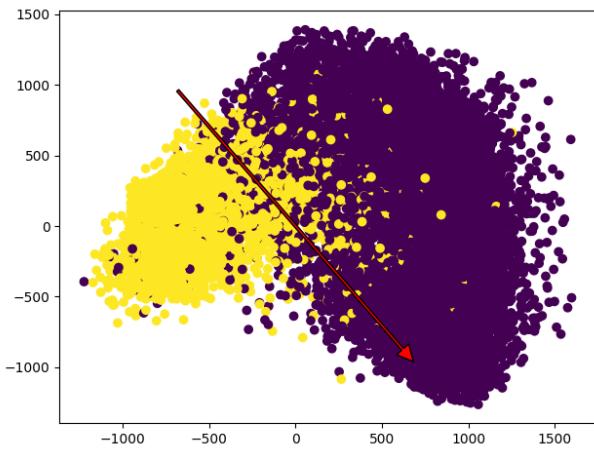


h) we expect that LDA projects data to the line which makes the most distinction between the two classes. As we see in the above figure, LDA has founded the projection line in the direction which most distinguishes the two classes. Actually it finds the line by maximizing the between class scatter and minimizing the within class scatter.

I) in the following plot, different classes are shown in the 2d space. As we see 10 different classes are not separated well and they are overlapping.



j, k) separator line is shown in the following figure:



and here are confusion matrixes:

train set	test set
<code>[[7513 502] [315 7723]]</code>	<code>[[1853 132] [98 1864]]</code>
0.94910	0.94172

l) most misclassified eyeglasses and pants are shown in the following table:

most misclassified eyeglass	most misclassified pant

m) required number of principal components for 90% variance is 170.

results of previous part with this number of principal components is below. As we see the performance of model is improved by using more principal components.

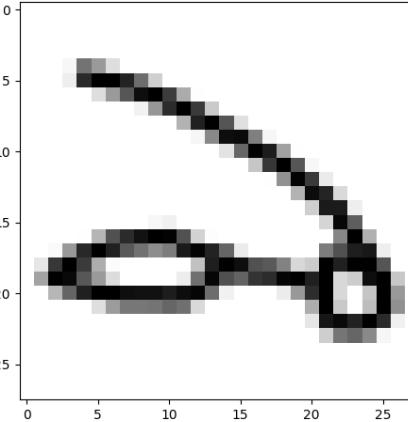
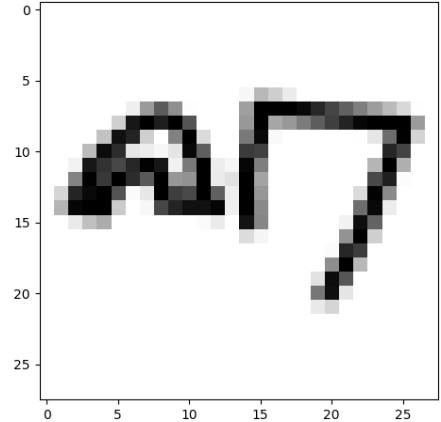
train set	test set
[[7791 224] [278 7760]]	[[1924 61] [76 1886]]
0.96872	0.9652900

N)

	1	3	5	7	9
train	[[7622 393] [405 7633]]	[[7724 291] [278 7760]]	[[7739 276] [261 7777]]	[[7746 269] [260 7778]]	[[7760 255] [260 7778]]
test	[[1885 100] [100 1862]]	[[1900 85] [79 1883]]	[[1905 80] [76 1886]]	[[1910 75] [80 1882]]	[[1911 74] [81 1881]]

According to the results, the best K for KNN algorithms is 7 because it leads to a minimum test error.

Misclassified eyeglasses and pants are shown in the following table:

misclassified eyeglass	misclassified pant
	

question 4)

a)

(a)

$$y = w_2 + \epsilon \cdot w_1 + b \cdot 1 = 9$$
$$w_1 = +.5, w_2 = -1$$

bias $w_3 = 1$

sign of y indicates the class.

activation = linear

b)

(b)

$$y = |x_2 - .5| + |x_2 - 1| - .5 \cdot 1 = 0$$
$$y = (x_2 - .5) \times (x_2 - 1) < 0 \rightarrow x_2^2 + .5x_2 - 1.1x_2 < 0$$

bias $w_3 = 1, w_2 = 1, w_1 = -.5$

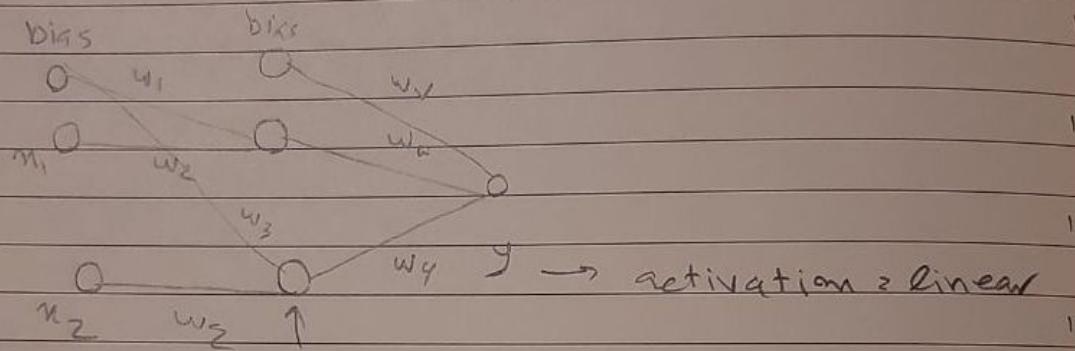
sign of y indicates the class.

activation = linear

$$b) y = \text{sign}(x_2 - w^T) + \text{sign}(\cdot, 1 - w_2) = 1$$

if $y = 1 \rightarrow \text{black}$

if $y < 0 \rightarrow \text{white}$



activation func = sign

$$w_1 = -1, w_2 = 1$$

$$w_3 = -1, w_4 = 1$$

$$w_0 = w_4 = 1, w_V = -1$$

c)

③ $m_1 > \omega$ and $m_2 < \nu$, $m_2 > 0$, $m_1 < 1$

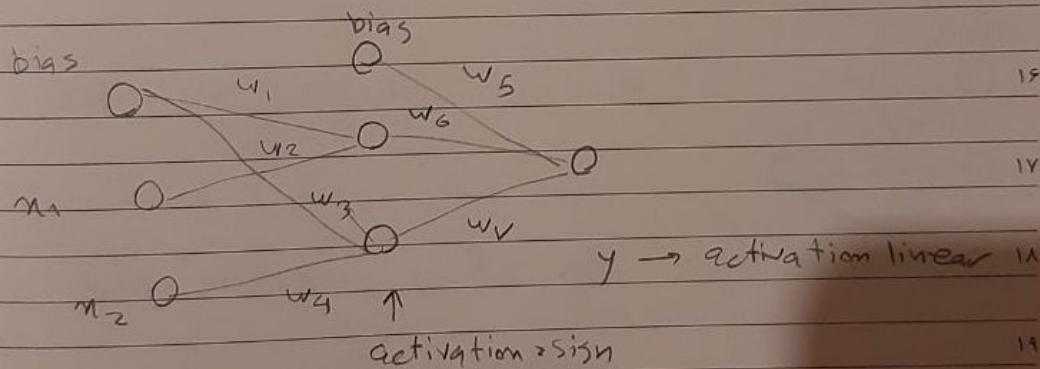
$(m_1 - \omega) > 0$ and $(\nu - m_2) > 0$

$$|m_1 - \omega| + |m_1 - 1| + |m_2| + |\nu - m_2| = 1, \delta = 0$$

$$y = \text{sign}(\omega - m_1) + \text{sign}(\nu - m_2) - 1$$

$y > 0 \rightarrow \text{black}$

$y \leq 0 \rightarrow \text{white}$

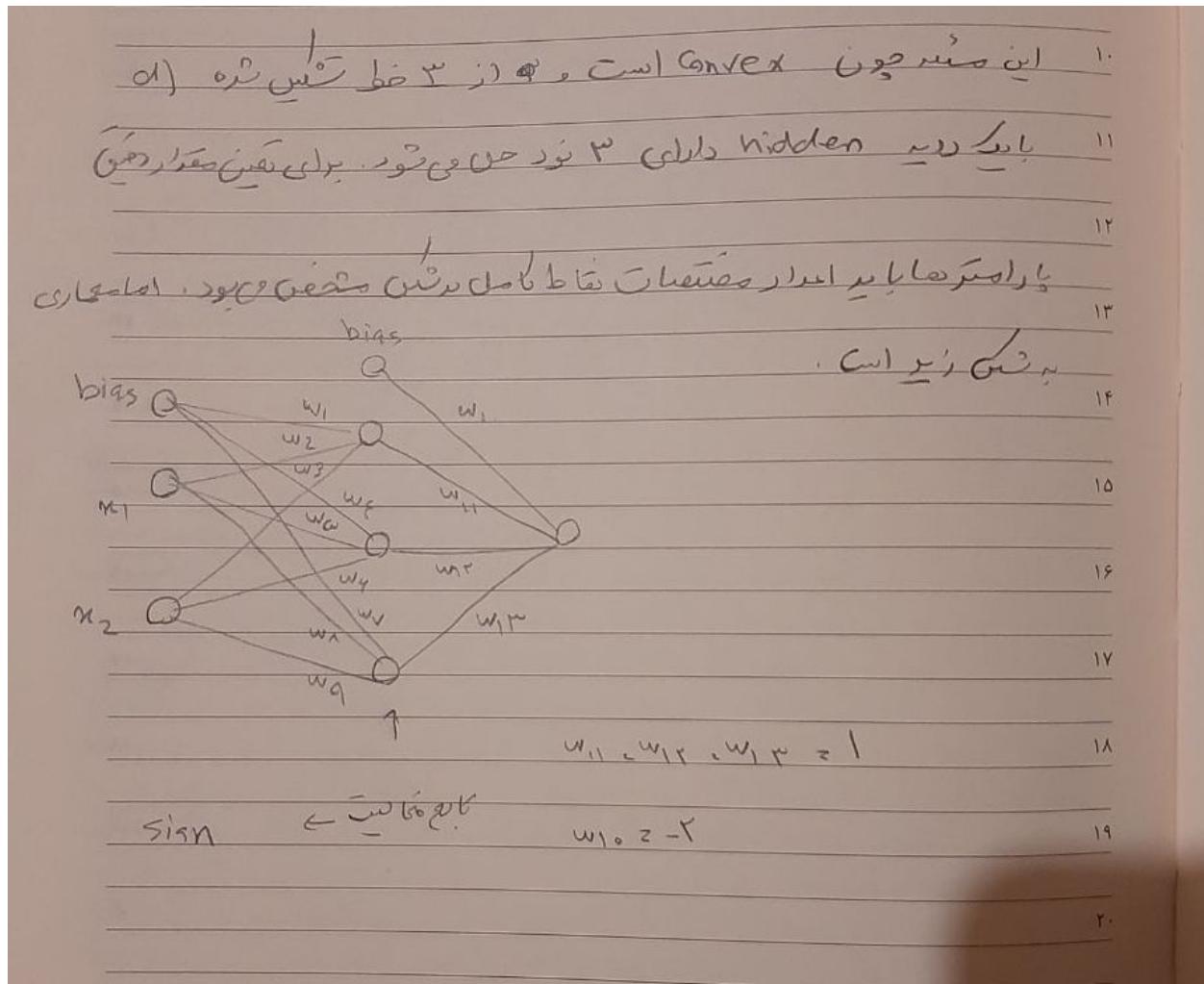


$$w_1 = -\omega, w_2 = 1, w_3 = \nu, w_4 = -1$$

$$w_5 = -1, w_6 = w_7 = 1$$

C	H	E	S	O	R	S	G	O
T	I					T		
A	V	P	S	T				
Y	10	14	13	11	12	11	10	9

d)



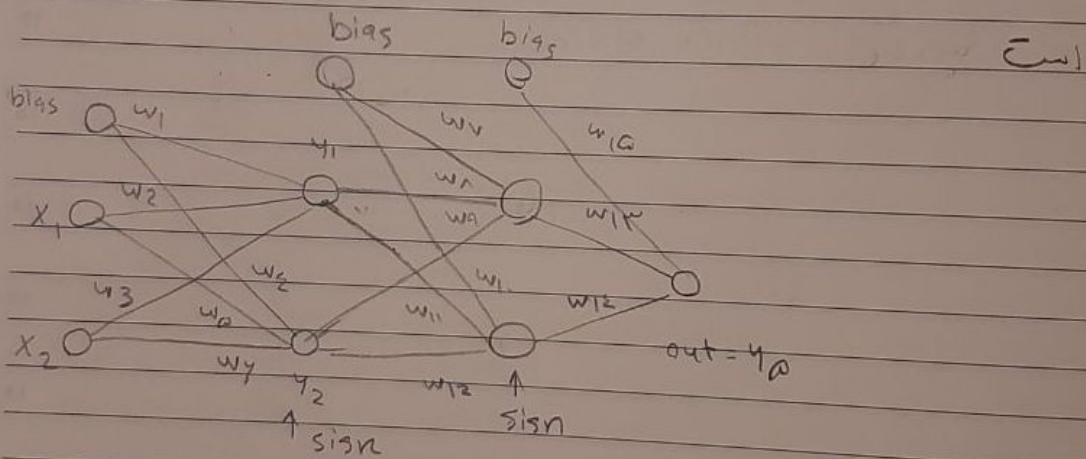
e)

$$\textcircled{1} \quad \begin{aligned} x_2 - yx_1 - 14 &\rightarrow x_2 - rx_1 + 14 = y_1 \\ x_2 = -rx_1 + 115 &\rightarrow x_2 + rx_1 - 115 = y_2 \end{aligned}$$

$$y_1 > 0, y_2 > 0 \rightarrow y_1 = \text{sign}(y_1) + \sin(\pi y_1) = 2 \\ y_2 = \text{sign}(y_2) + \sin(\pi y_2) = 1$$

$$\begin{aligned} y_1 < 0, y_2 < 0 \quad \rightarrow \quad & \text{sign}(y_1) + \text{sign}(y_2) = -2 \\ & \text{sign}(-y_1) + \text{sign}(-y_2) = 2 \\ & y_2 = \text{sign}(-y_1) + \text{sign}(-y_2) - 1 \end{aligned}$$

در ماتریس هم راز دو حالت با ۱۷ معناد



$$\begin{array}{lll} w_1 = 4 & w_{22} = 1, 3 & w_V = w_{1,2} - 1 \\ w_2 = -4 & w_{02} = 4 & w_{VW} = w_9 = 1 \\ w_3 = 1 & w_4 = 1 & w_{11} = w_{12} = 1 \end{array}$$

$$w_1 \omega, w_1 \nu, w_1 \sigma \geq 1$$

ش	ی	د	س	ج	ب	غ	ق
۲۱							
۲	۱						
۳	۴	۵	۶	۷	۸	۹	
۱۰	۱۱	۱۲	۱۳	۱۴	۱۵	۱۶	
۱۷	۱۸	۱۹	۲۰	۲۱	۲۲	۲۳	
۱۸	۱۹	۲۰	۲۱	۲۲	۲۳	۲۴	
۲۴	۲۵	۲۶	۲۷	۲۸	۲۹	۳۰	

$y_0 = 1 \rightarrow$ belongs to black
 $y_0 = 0 \rightarrow$ belongs to white

- f) this network should at least has five following inputs: body height, body length, eye diameter, horn length and leg length. We can also extract better features from these features and give it as input to the network.
- g) As this problem is a classification problem and there are 5 different classes in this problem, this network needs to have 5 outputs where each one indicates whether the given data belongs to that class or not.
- h) Yes, this network should have at least one hidden network. Optimum count of hidden units would be determined by experiment. Actually deeper networks solve the problem better but it also depends on the size of the data. For training deep networks we need more data! We should also take into account over-fitting phenomena when determining the best depth.
- i) Recurrent connections are useful when we have a series data, like a sequence of words, prices and etc. in this case we don't need recurrent connection and a deep feed forward network solves the problem.
- j) Different activation functions like ReLu, logistic, TanH, sigmoid and softmax can be used for this problem, but usually softMax is used for multi-classification problem.
<https://towardsdatascience.com/activation-functions-neural-networks-1cbd9f8d91d6>
- k) Different learning algorithms like Adam, Adagrad, RMSProp, and SGD can be used for this problem.

question 5)

a)

```
38      v=b;
39      for j=1:nSV
40          v = v + alpha(j) * kernel_fnc(SV(j,:)',xy)
41      end
42
43      if ~v < 0
```

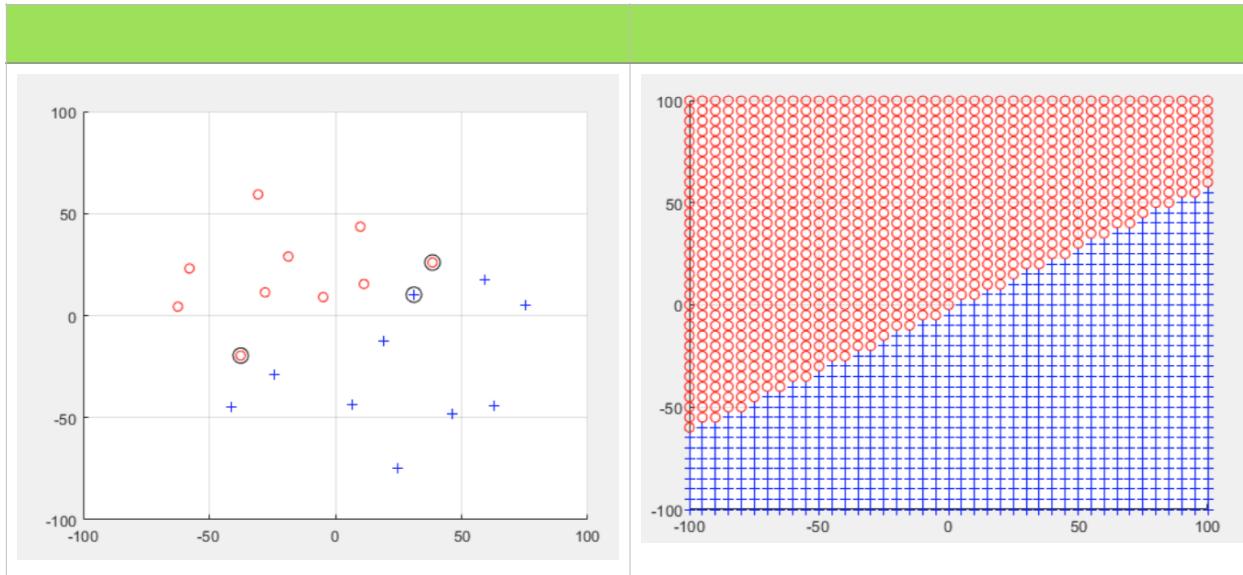
b)

```
42      for j=1:n
43
44          xj = D(j,1:2);
45          yj = D(j,3);
46
47          H(i,j) = kernel_fnc(xi',xj') * yj * yi ;
48
49      end
```

c)

```
75 options = optimset('LargeScale','off');
76
77 [x,fval] = quadprog(H, f, A, b, Aeq, beq, lb, ub, x0, options);
78
79 tol = 1e-8;
```

d)



outputs of this function are described below:

- 1) SV: is array of support vectors
- 2) alpha : is array of parameter alphas in SVM algorithm
- 3) b : is also parameter b in SVM

e) circles points indicate the support vectors.

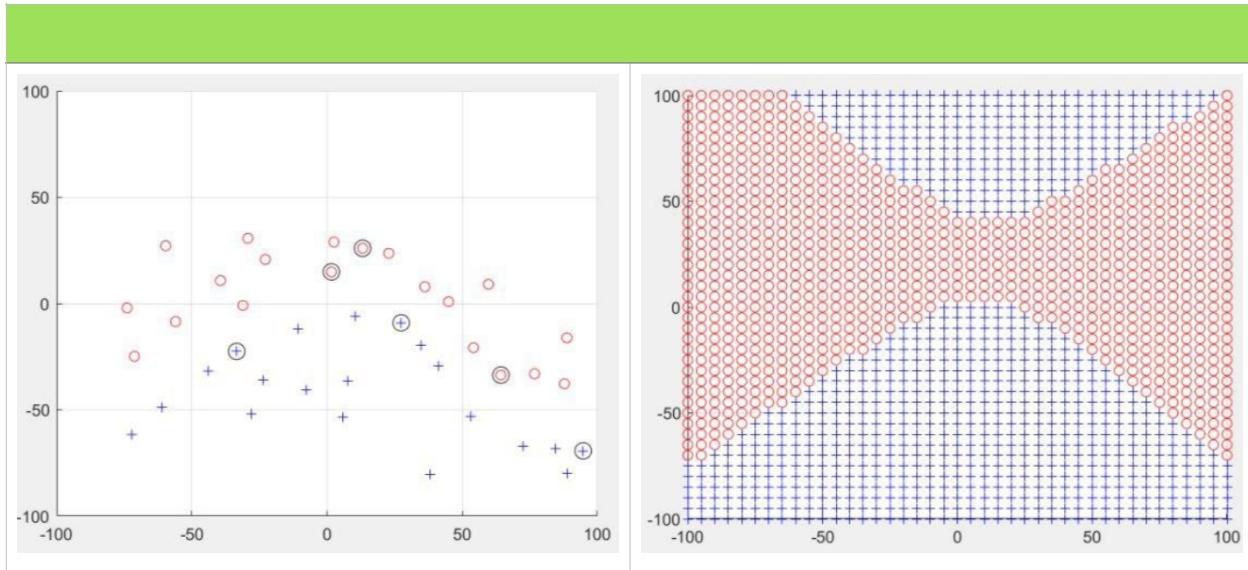
f) according to the comments in code, inf indicates cost of misclassifying. It is used for determining the cost of slack variable.

```
14 % Format x_1, x_2, y per line
15 %
16 % C is the slack cost -- the cost of misclassifying training point in
17 % the non-separable case. C = inf means separable case.
18
19
```

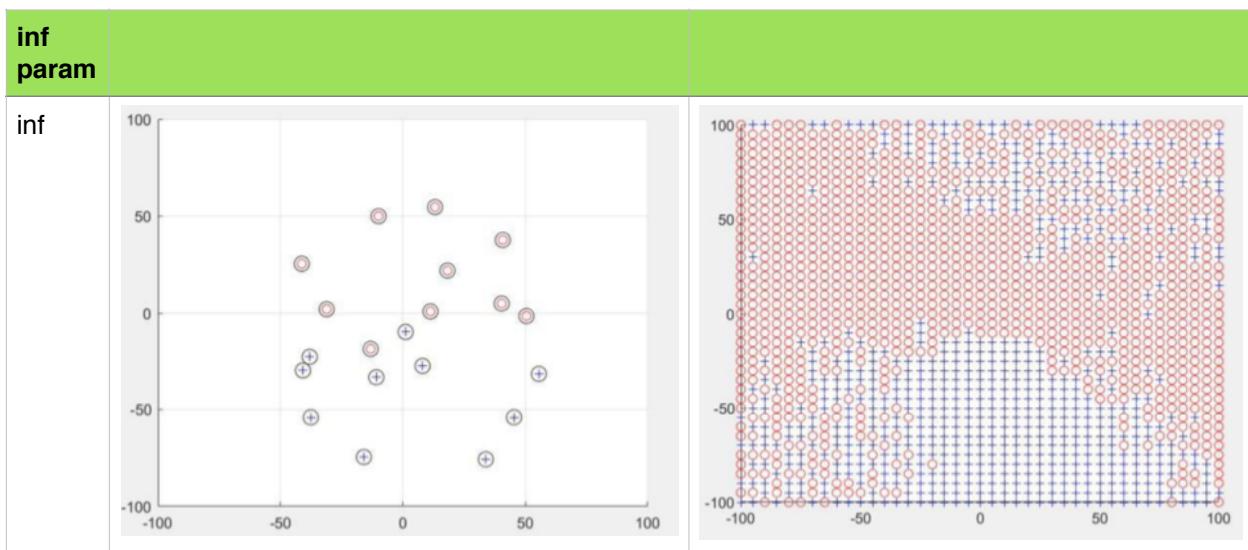
g) this new kernel function is as below:

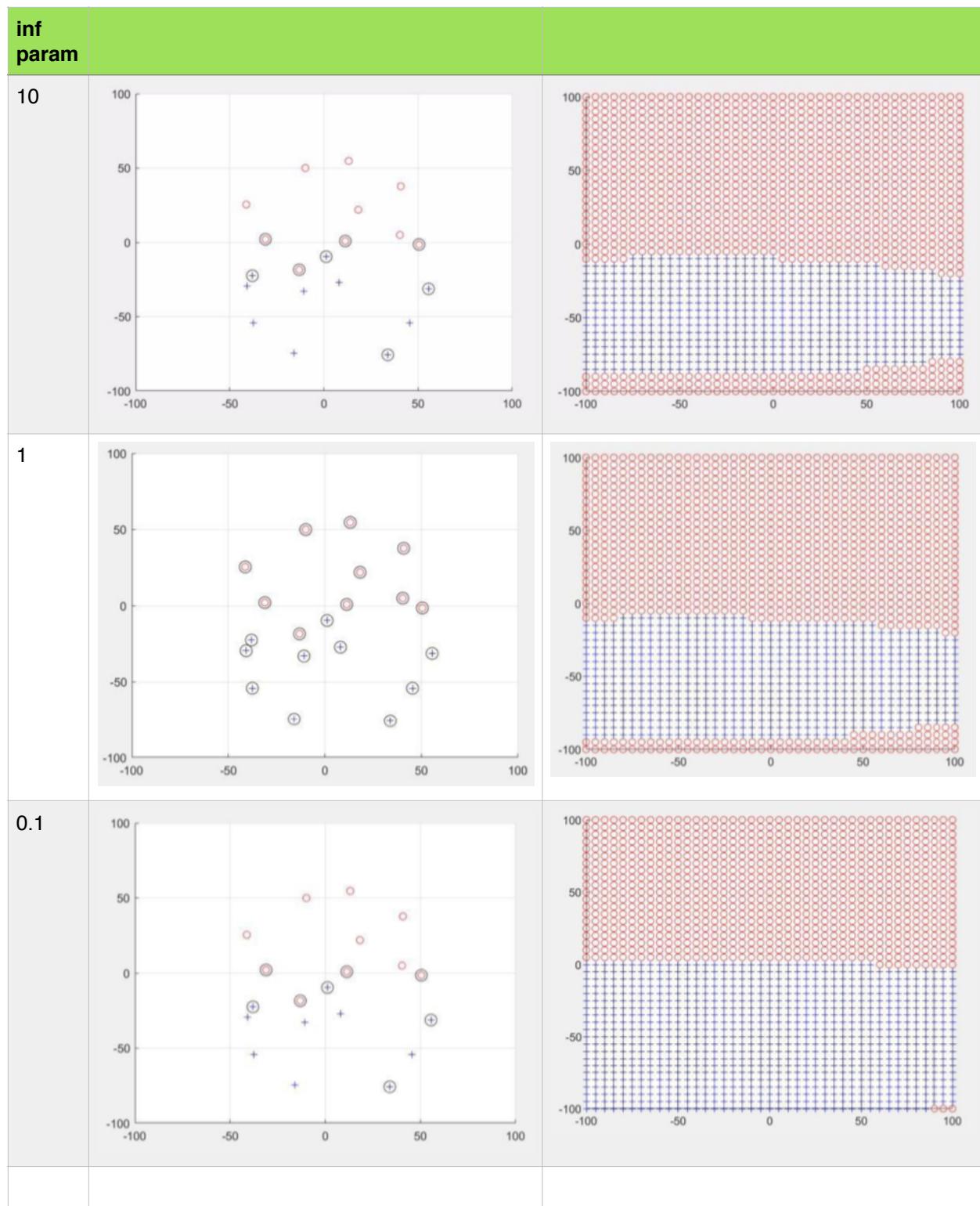
```
1 function p = quad_kernel(x,y,args)
2
3 p = (x'*y + 1)^2;
4
```

result is shown in the following:



h)





C is a regularization parameter that controls the trade off between the achieving a low training error and a low testing error that is the ability to generalize the classifier to unseen data.

Consider the objective function of a linear SVM : $\min \|w\|^2 + C \sum \xi_i$. If C is too large the optimization algorithm will try to reduce $\|w\|$ as much as possible leading to a hyperplane which tries to classify each training example correctly (similar to first figure). Doing this will lead to loss in generalization properties of the classifier. On the other hand if C is too small then you give your objective function a certain freedom to increase $\|w\|$ a lot, which will lead to large training error.

C Parameter is used for controlling the outliers, low C implies we are allowing more outliers, high C implies we are allowing fewer outliers.

question 6)

A)

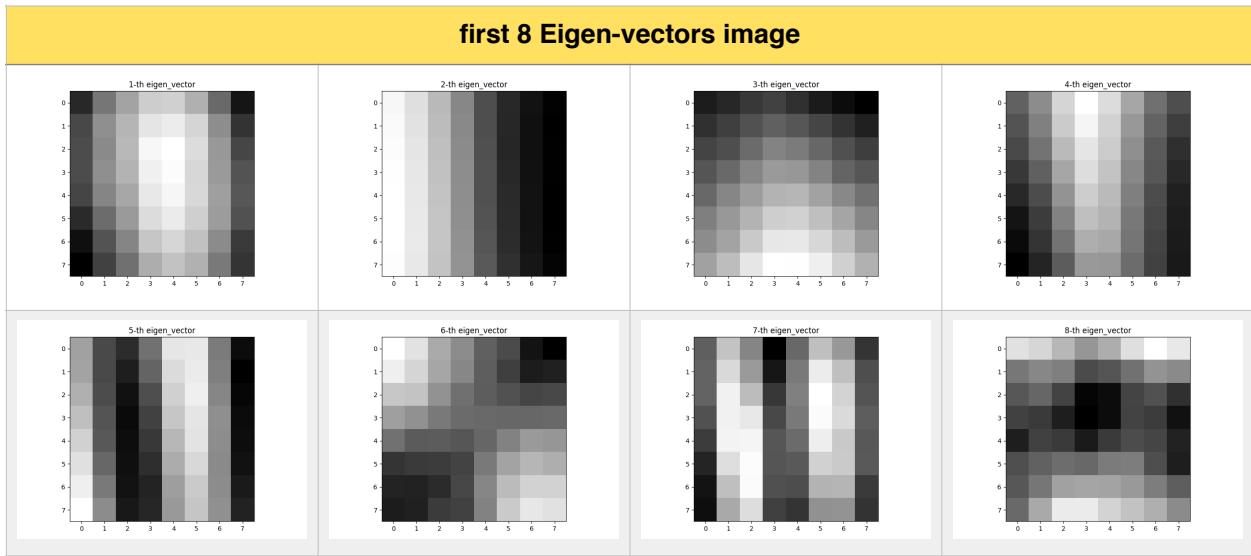
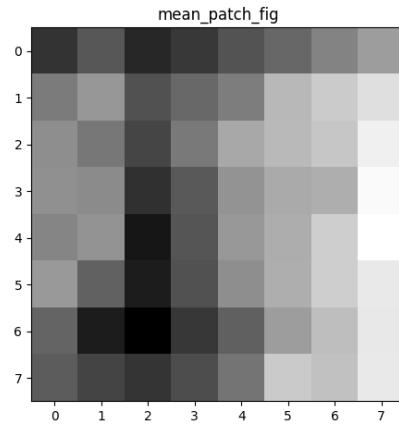
- B) first 20 largest eigenvalues and corresponding eigenvectors are written in the 'eigenval_eigenvec.txt' file:

```
eig val: 327601.15508234745
eig vec: [ 1.21470474e-01  1.73473707e-01 -1.97857292e-01 -4.34302750e-02
 6.70055789e-02  2.72696143e-01 -7.70942106e-02  2.20365301e-01
-2.65718351e-01 -9.95971150e-02 -9.33477740e-02  3.61785507e-01
-4.60718447e-02  1.94580802e-01  1.62961170e-01 -3.54577586e-01
-1.57810692e-01 -1.48831563e-01  2.36592863e-02 -1.19570534e-01
-3.04814360e-01  6.53660587e-03  3.93008709e-02  4.86440969e-02
-6.17373402e-03 -1.49178032e-01  7.12543990e-02 -1.78302439e-01
-5.99033999e-02 -1.14680527e-01 -2.70295754e-01 -7.09294952e-02
 1.18994382e-01 -1.44437863e-02 -2.20179847e-02  1.01787393e-01
 3.49875508e-02 -4.92339033e-02  8.29504280e-03 -1.29284066e-02
-5.02296792e-03 -1.06887453e-01  5.71634550e-03  2.04938432e-02
-3.72728299e-06  1.42092492e-02 -5.72518187e-03  2.87987221e-02
 1.77199402e-02  3.69338714e-02 -4.65175182e-02 -5.75438864e-03
 1.41029383e-03  5.64168644e-03 -2.53447298e-03  2.77942939e-02
 6.17545874e-02  1.98244065e-03  3.04468482e-02 -4.62164885e-02
-1.18169617e-02 -2.13719427e-02  3.67555311e-02  4.09907244e-03]

eig val: 17556.24824832425
eig vec: [ 0.12400416  0.14040671 -0.17597622  0.03910798 -0.07241145  0.21411328
 0.09246566  0.20004021 -0.0390126 -0.0193197 -0.21794224  0.11533818
-0.06019659 -0.00459239 -0.08804577 -0.14616581 -0.14719846  0.16949958
-0.0450442 -0.04912508  0.4044589 -0.00535826  0.02160375 -0.15487313
 0.03456892 -0.20372945 -0.07694805 -0.04107012  0.07059579  0.07089179
 0.26533025  0.34772974 -0.13837784  0.10816922 -0.14071918 -0.16719343
-0.22928711  0.00888375 -0.01582016 -0.0071815  0.01990884  0.04818369
 0.060081 -0.05429138 -0.09300315  0.01649849  0.05681312 -0.11029378
-0.07005266 -0.13589187  0.09703994 -0.038294  0.0047069 -0.03263127
 0.05978006 -0.02403855 -0.02910692  0.0316934 -0.08347474  0.02780222
 0.02067231  0.12603213 -0.09014985  0.06390433]

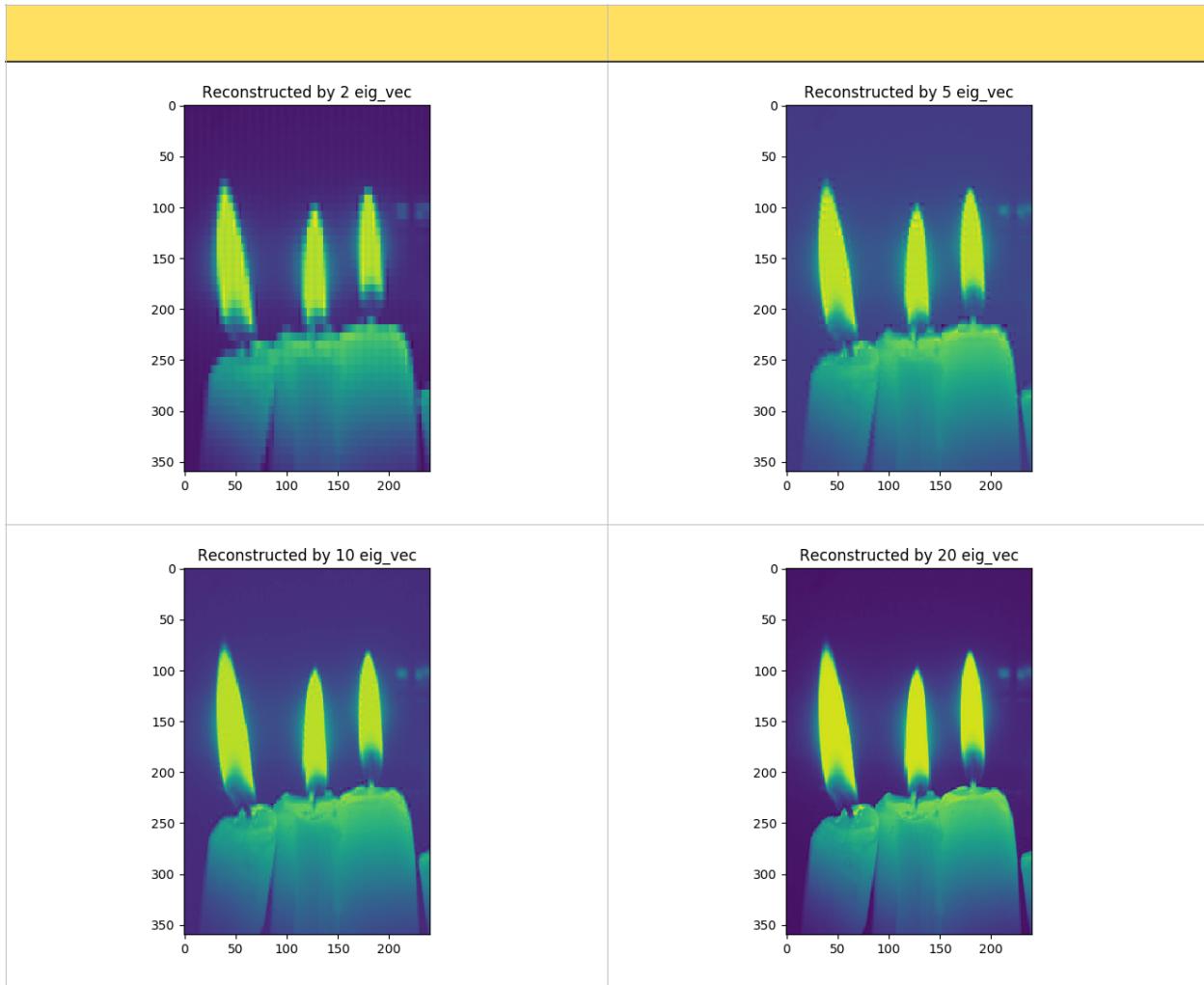
eig val: 3141.5701634757343
eig vec: [ 0.1254903  0.08832029 -0.14119957  0.17777289 -0.11734737  0.10314953
-0.01326631  0.13937874  0.19565377 -0.17622927 -0.19850112 -0.20629534
-0.12220621 -0.08068153 -0.07242062 -0.07480381 -0.13810181  0.14337822
```

Here is the image of mean patch:



C, D, E)

Result is shown in the following table:

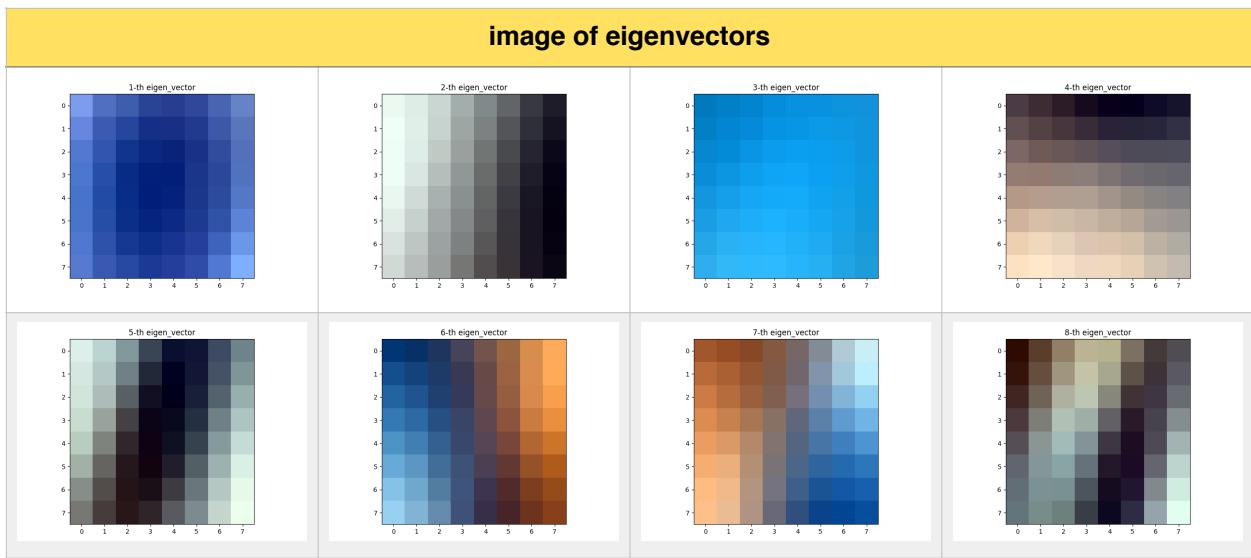
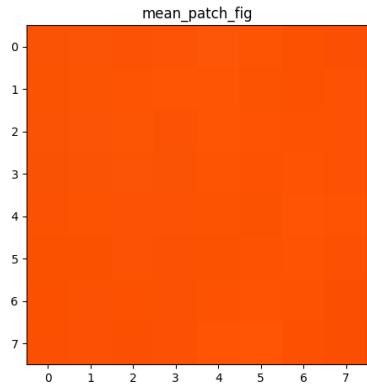


As we see by increasing the count of eigen-vectors quality of reconstructed image increases. Actually quality of the image is related to the variance ratio of used eigenvectors not number of used eigenvectors. based non the above figures quality doesn't change more by adding 10 eigenvectors after first 10 eigenvectors.

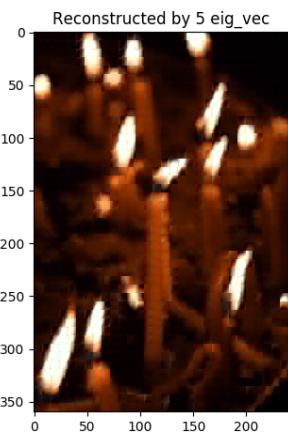
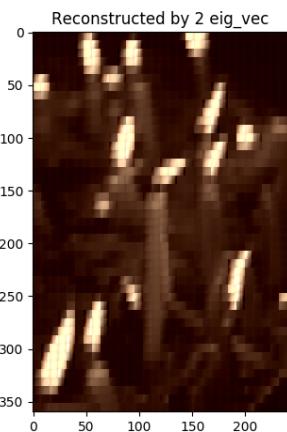
F)

eigenvalues and eigenvectors of this question are written in file (eigenval_eigenvec_f.txt).

patch image:



Reconstructed images



question 7)

- A. No it doesn't make sense, Because PCA is used for dimensionality reduction not for increase in dimensionality. also it works by extracting the eigenvectors and eigenvalues and it is not possible to extract Eigen vectors more than the size of dimensionality or size of the dataset.
- B. consider two vectors W_1 and W_2 . $\text{span}(W_1, W_2)$ is $z_1W_1 + z_2W_2$ for scalars z_1 and z_2 . in PCA we have label switching which means $\text{span}(W_1, W_2) = \text{span}(W_2, W_1)$. Then PCA has no uniqueness because of the label switching which means we can exchange rows of matrix W with each other. changing rows of W means changing the order of factors which is possible. (<https://www.cs.ubc.ca/~schmidtm/Courses/340-F17/L27.pdf>) (W is the transformation matrix)
- c) SVD converts the main matrix to three matrices like USV , but PCA eliminates less significant components and transforms data into a new space where properties are not correlated there. Relationship between SVD and PCA is described below. (Ref: <https://intoli.com/blog/pca-and-svd/>)

Let the data matrix X be of $n \times p$ size, where n is the number of samples and p is the number of variables. Let us assume that it is centered, i.e. column means have been subtracted and are now equal to zero.

Then the $p \times p$ covariance matrix C is given by $C = X^T X / (n - 1)$. It is a symmetric matrix and so it can be diagonalized:

$$C = V L V^T$$

where V is a matrix of eigenvectors (each column is an eigenvector) and L is a diagonal matrix with eigenvalues λ_i in the decreasing order on the diagonal. The eigenvectors are called *principal axes* or *principal directions* of the data. Projections of the data on the principal axes are called *principal components*, also known as *PC scores*; these can be seen as new, transformed, variables. The j -th principal component is given by j -th column of XV . The coordinates of the i -th data point in the new PC space are given by the i -th row of XV . If we now perform singular value decomposition of X , we obtain a decomposition

$$X = USV^T$$

where U is a unitary matrix and S is the diagonal matrix of singular values s_i . From here one can easily see that:

$$C = V \frac{S^2}{n-1} V^T$$

meaning that right singular vectors V are principal directions and that singular values are related to the eigenvalues of covariance matrix via $\lambda_i = s_i^2 / (n-1)$. Principal components are given by $XV = USV^T V = US$.

To summarize:

1. If $X = USV^T$, then columns of V are principal directions/axes.
2. Columns of US are principal components ("scores").
3. Singular values are related to the eigenvalues of covariance matrix via $\lambda_i = s_i^2 / (n-1)$. Eigenvalues λ_i show variances of the respective PCs.
4. To reduce the dimensionality of the data from p to $k < p$. select k first columns of U , and $k \times k$ upper-left part of S . Their product $U_k S_k$ is the required $n \times k$ matrix containing first k PCs.

d) We know that PCA ignores small values of Eigen-values which leads to numerical instability. As singular values are more numerical stable than eigenvalues, in this case it would be better to use SVD for performing PCA instead of calculating the XX and performing Eigen-value and Eigen-vector decomposition. In the previous part it is explained completely how to use SVD for performing PCA. But here is a summary:

SVD on matrix X results in obtaining matrices U, S, and V as:

$$X = USV$$

where columns of US are principal components.

using SVD may take longer time but it has higher numerical accuracy.

e) PCA can reduce noise but it can't eliminate the noise. In PCA an orthogonal linear transformation is used to find a projection of all data into k dimensions, whereas these k dimensions are those of the highest variance. The eigenvectors of the covariance matrix (of the dataset) are the target dimensions and they can be ranked according to their eigenvalues. A high eigenvalue signifies high variance explained by the associated eigenvector dimension. This process of keeping large eigenvalues and corresponding eigenvectors helps to reduce the noise and keep the most meaningful information of image.

f) Back Propagation is used to effectively train a neural network through a method called chain rule. In simple terms, after each forward pass through a network, back propagation performs a backward pass while adjusting the model's parameters (weights and biases). If the used activation function in network is not differentiable then the BP is not applicable.

Networks like Restricted Boltzmann Machine (RBM) and Hard attention can't use BP.

g) Yes it is possible. It is differentiable and has returns a value between 0 and 1. But it is not used in deep neural networks because of following problems:

- 1) It kills gradients at saturation point, and cause vanishing gradient problem,
- 2) Its outputs are not zero-centered.

h) Yes it can be used for regression. SVM works based on maximizing the margin. Maximizing the margin can more generally be seen as regularizing the solution by minimizing W (which is essentially minimizing model complexity) this is done both in the classification and regression. But in the case of classification this minimization is done under the condition that all examples are classified correctly and in the case of

regression under the condition that the value Y of all examples deviates less than the required accuracy ϵ from $f(x)$ for regression.

In Regression the goal is to find a function $f(x) = wx+b$ under the condition that $f(x)$ is within a required accuracy ϵ from the value $y(x)$ of every data point, i.e. $|y(x) - f(x)| \leq \epsilon$ where ϵ is the distance between the red and the grey line. Under this condition we again want to minimize $f'(x)=w$, again for the reason of regularization and to obtain a unique solution as the result of the convex optimization problem. One can see how minimizing w results in a more general case as the extreme value of $w=0$ would mean no functional relation at all which is the most general result one can obtain from the data.

