

Theory Questions

(1) As suggested we will choose

$$m = d = \left\lfloor \frac{1}{\gamma^2} \right\rfloor \quad \text{and } \{x_i\}_i \text{ to be}$$

The standard basis of \mathbb{R}^d , namely $\{\hat{e}_i\}_i$

e.g. $\hat{e}_1 = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$ Then:

(2) by definition $|x_i| = |\hat{e}_i| = \gamma$

(3) we would like to find w^* s.t
 $y_i x_i w^* \geq \gamma$ and $|w^*| \leq 1$.

$$* | \sum_i w_i^* x_i |^2 = d w_i^* \Rightarrow w_i^{*2} \leq \frac{1}{d} = \frac{1}{\left\lfloor \frac{1}{\gamma^2} \right\rfloor}$$

assume for simplicity $w_i = w_j \forall i, j$ see next page

* also we know from (c) $(x_i, y_i) = (\hat{e}_i, -1)$

$$\Rightarrow \frac{y_i x_i w^*}{|w^*|} = -\frac{w_i^*}{|w^*|} \geq \gamma \Rightarrow w_i^* \leq -\frac{\gamma}{|w^*|}$$

$$|w^*| \geq 0, 0 < \gamma < 1 \Rightarrow w_i^* \leq 0$$

$$* + ** \Rightarrow w_i^* \leq -\sqrt{\frac{1}{\left\lfloor \frac{1}{\gamma^2} \right\rfloor}}$$

we choose

$$\Rightarrow w^* = \frac{1}{\sqrt{\left\lfloor \frac{1}{\gamma^2} \right\rfloor}} (-1, -1, \dots, -1)$$

(c) we have to choose $\{y_i\}$; so that the perceptron will make $\lfloor \frac{1}{\gamma^2} \rfloor$ mistakes
 $\lfloor \frac{1}{\gamma^2} \rfloor = m = d \Rightarrow$ mistakes for all x_i !

we start with $w_1 = (0, \dots, 0)$.

$$\text{sign}(w_1 \cdot x_1) = 0 \Rightarrow \hat{y}_1 = 1$$

so we can choose $y_1 = -1$ for the perceptron to make a mistake.

\Rightarrow non-zero update:

$$w_2 = w_1 + y_1 x_1 = (-1, 0, \dots, 0)$$

we can repeat that logic and then we expect:

$$w_t = (-1, -1, \dots, -1, \underset{\uparrow}{0}, \dots, 0)$$

t index

We can prove it by induction. We already have shown for $t=2$. Now assuming our logic. We have w_t as described.
(and $y_t = -1$)

$$\text{sign}(w_t \cdot x_t) = \text{sign} [(-1, -1, \dots, -1, \underset{\uparrow}{0}, \dots, 0) \cdot (0, \dots, 0, \underset{\uparrow}{1}, 0, \dots, 0)] = 0$$

t entry t entry

$\Rightarrow \hat{y}_t = 1$, again we choose $y_t = -1$ for making another mistake and then

$$\begin{aligned} w_{t+1} &= w_t + y_t x_t = (-1, \dots, -1, \underset{\uparrow}{0}, \dots, 0) + (-1, \underset{\uparrow}{0}, \dots, 0, \underset{\uparrow}{1}, 0, \dots, 0) \\ &= (-1, \dots, -1, \underset{\uparrow}{0}, \dots, 0) \end{aligned}$$

t t'

\Rightarrow our point are $\{(\hat{x}_i, -1)\}_{i=1}^m$, and we have
 $m = \lfloor \frac{1}{\gamma^2} \rfloor$ mistakes

$$(2) \quad d \geq 6 \quad x \in \mathbb{Z} = \{1, 2, \dots, d\} \quad H = \{h_{ij} \mid 1 \leq i < j \leq d\}$$

$$L_{ij}(x) = \begin{cases} 1 & (x=i) \cup (x=j) \\ 0 & \text{otherwise} \end{cases}$$

define $n_1 = |\{h \in H \mid L(x) = 1\}|$
 $n_0 = |\{h \in H \mid L(x) = 0\}|$

we start with $|H_1| = \binom{d}{2} = \frac{d(d-1)}{2}$

assuming we got a point x_1 ,

then $|n_1| = |\{h_{x_1, x} \mid x > x_1\} \cup \{h_{x < x_1, x}\}| = d - 1$

and $|n_0| = |H_1| - |n_1| = \frac{d(d-1)}{2} - (d-1)$

$$= \left(\frac{d}{2} - 1\right)(d-1) \geq 2(d-1) > |n_1|$$

\uparrow
 $d \geq 6$

\Rightarrow first choice of label is $\hat{y}_1 = 0$

assuming we made a mistake and $y_0 = 1$
 $\Rightarrow |H_2| = |H_1| - |n_0| = |n_1| = d - 1$

$$H_2 = \{h_{x_1, x} \mid x > x_1\} \cup \{h_{x < x_1, x}\}$$

now assuming we have a new point $x_2 \neq x_1$,
 there is only one hypothesis for which $\hat{y}_2 = 1$

$$\Rightarrow |n_1| = 1, |n_0| = d - 1 - |n_1| = d - 2$$

$$\Rightarrow |n_0| > |n_1| \Rightarrow \hat{y}_2 = 0.$$

\uparrow
 $d \geq 6$

but if we made a mistake and $y_2 = 1$

$$\Rightarrow |H_3| = |H_2| - |n_0| = 1 \Rightarrow \text{we are done!}$$

\Rightarrow we showed that at most, we will have 2 mistakes until we are left with only one hypothesis. If at any stage we wouldn't make a mistake we will be left with even smaller hypothesis class size. So the case of 2 mistakes in a row is the most striking.

Also, we can think of an example for such a sequence. (its labels shall be =1) so $(3, 1), (2, 1)$ will do.

After the first we are left with h_{13}, h_{23} and $\{h_{3x}\}_{x \geq 3}$ and for the second we only filter h_{34} .

$\Rightarrow \text{Sup}_{\mathcal{S}} M_{\Delta}(\mathcal{S}) = 2$ ~~/~~

(3) f_1, \dots, f_m are convex $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$

$$g(x) = \max_i f_i(x)$$

(a) define $j = \arg \max_{i \in [1, m]} f_i(\alpha x_1 + (1-\alpha)x_2)$
 $0 < \alpha < 1$

$$\Rightarrow g(\alpha x_1 + (1-\alpha)x_2) = f_j(\alpha x_1 + (1-\alpha)x_2)$$

$$\leq \alpha f_j(x_1) + (1-\alpha)f_j(x_2) \leq \alpha g(x_1) + (1-\alpha)g(x_2)$$

f_j is convex

$$\downarrow$$

$$f_j(x) \leq g(x) \quad \forall x$$

$\Rightarrow g(x)$ is convex. //

(b) sub gradient set $\{\nabla f_i\}$ (at point x_1)

$$\forall x \quad f(x) \geq f(\tilde{x}) + \nabla f(\tilde{x})(x - \tilde{x})$$

$$g(x) \geq f_j(x) \geq f_j(\tilde{x}) + \nabla f_j(\tilde{x})(x - \tilde{x}) =$$

\downarrow
 f_j is convex

$$= g(\tilde{x}) + \nabla f_j(\tilde{x})(x - \tilde{x})$$

$$j = \arg \max_{j \in [1, m]} \{f_j(\tilde{x})\}_{j=1}^m$$

$$\min_{w,b,g} \frac{1}{2} w^T w + \frac{C}{2} \sum_{i=1}^m g_i^2$$

$$\text{s.t. } y_i(w^T x_i + b) \geq 1 - g_i \quad \forall i = 1 \dots m$$

(a) we have to show that some other problem, similar to the original but with the additional set of constraints $g_i \geq 0$ is equivalent.

Assume (in contradiction) that there is an optimal solution to the original problem in which there exist some k for which $g_k < 0$.

we can think of some other set \tilde{g}_i such that $\tilde{g}_j = g_j \quad \forall j \neq k$ and $\tilde{g}_k = 0$. This set of variables also satisfies the constraint

$$y_j(w^T x_j + b) \geq 1 - g_j \quad \tilde{g}_j = 1 - g_j \quad j \neq k \\ \geq 1 - \tilde{g}_j \quad j = k$$

But it has a smaller optimum because the element $\frac{C}{2} g_k^2$ now vanish

and for the original set $\frac{C}{2} g_k^2 > 0$

\Rightarrow contradiction to the assumption that $\{g_i\}$ are of the optimal solution

$\Rightarrow g_i \geq 0 \quad \forall i$ must hold anyway!

⑥ The Lagrangian:

$$f(x_i, y_i, w, b, \xi_i, \lambda_i) =$$

$$= \frac{1}{2} w^T w + \frac{C}{2} \sum_i \xi_i^2 + \sum_i \lambda_i (1 - \xi_i - y_i (w^T x_i + b))$$

$$\textcircled{1} \quad \nabla_w f = w - \sum_i \lambda_i y_i x_i \stackrel{!}{=} 0 \quad \left. \begin{array}{l} \text{minimizing} \\ \text{w.r.t.} \end{array} \right\}$$

$$\Rightarrow w = \sum_i \lambda_i y_i x_i$$

$$\nabla_{\xi_i} f = C \xi_i - \lambda_i \stackrel{!}{=} 0$$

$$\Rightarrow \xi_i = \frac{\lambda_i}{C}$$

$$\nabla_b f = - \sum_i \lambda_i y_i \stackrel{!}{=} 0$$

$$\Rightarrow \sum_i \lambda_i y_i = 0$$

for optimal value

$$\Rightarrow f = \frac{1}{2} \underbrace{\left(\sum_i \lambda_i y_i x_i \right)}_{=0} \underbrace{\left(\sum_j \lambda_j y_j x_j \right)}_{\lambda_j \cdot \xi_j} + \frac{C}{2} \sum_i \lambda_i^2$$

$$= b \underbrace{\sum_i \lambda_i y_i}_{=0} + \underbrace{\sum_i \lambda_i}_{\lambda_i \cdot \xi_i} - \underbrace{\sum_i \lambda_i^2}_{\lambda_i^2} - \underbrace{\sum_i \lambda_i y_i x_i}_{w^T} \underbrace{\sum_j \lambda_j y_j x_j}_{w^T}$$

$$= -\frac{1}{2} \sum_i \lambda_i y_i x_i \sum_j \lambda_j y_j x_j - \frac{1}{2C} \sum_i \lambda_i^2 + \sum_i \lambda_i$$

d)

Our primal problem is

$$g(\lambda) = \min_{w, b, g} L(w, b, g, \lambda) \quad (\text{as seen in class})$$

where $g(\lambda)$ is the dual function.

The dual problem then is

$$\max_{\lambda \geq 0} g(\lambda) = \min_{\lambda \geq 0} -g(\lambda) \quad *$$

This is because the problem is

$$\min_{w, b, g} \max_{\lambda} L(w, b, g, \lambda)$$

* with the constraint from minimization of L

$$\sum \lambda_i y_i = 0 \quad \left(\begin{array}{l} \lambda \geq 0 \text{ to satisfy } \\ g_i \geq 0 \text{ by } \lambda_i = c g_i \end{array} \right)$$

\Rightarrow

$$\min_{\lambda \geq 0} \frac{1}{2} \sum_{ij} \lambda_i y_i x_i \lambda_j y_j x_j + \frac{1}{2C} \sum_i \lambda_i^2 - \sum_i \lambda_i$$

$$\text{s.t. } \sum \lambda_i y_i = 0$$

of course f , and g also depend on x_i, y_i . The inputs of the original function

$$(5) \quad H = \{h_{a,b} : a < b\} \quad X = \mathbb{R}$$

$$h(x) = \begin{cases} 1 & x \in [a, b] \\ 0 & \text{else} \end{cases}$$

$$J_H(m) = \max_{C: |C|=m} |H_C|$$

for convenience, we can get a sorted sample of m points. (x_1, x_2, \dots, x_m)

now, we can think of a possible labeling
 $\star (0, 0, \dots, 0) \rightarrow x_i \in [a, b] \quad \forall i.$

$\star \star$ other possibilities:

we can start from $x_1 \in [a, b]$,

There are m possibilities to label like that. $[a, b]$ can include points from x_1 up to some x_K ($K = [1, m]$).

for $x_2 \in [a, b]$ and $x_1 \notin [a, b]$ we have $m-1$ possibilities. points up to x_K ($K = [2, m]$) can be included. and so on. This is an algebraic series. \Rightarrow

$$J_H(m) = 1 + \sum_{k=1}^m (m+1-k) = 1 + \frac{m}{2}(m+1)$$

$\star \star$

⑥ we start with

$$n \geq \frac{256d}{\varepsilon^2} \ln\left(\frac{128d}{\varepsilon^2}\right) + \frac{256}{\varepsilon^2} \ln\left(\frac{1}{\delta}\right)$$

using $x \geq a \ln(x) + b$
 $\Rightarrow x > a \ln(x) + b$

we get $(a = \frac{64d}{\varepsilon^2}, b = \frac{128}{\varepsilon^2} \ln\left(\frac{1}{\delta}\right))$

$$n \geq \frac{64d}{\varepsilon^2} \ln(n) + \frac{128}{\varepsilon^2} \ln\left(\frac{1}{\delta}\right)$$

$$(\delta < 1) \geq \frac{64d}{\varepsilon^2} \ln(n) + \frac{32}{\varepsilon^2} \ln\frac{1}{\delta} \geq$$

$$\left| \begin{array}{l} d \geq \\ \varepsilon < 1 \end{array} \right\} \geq \frac{64d}{\varepsilon^2} \ln(n) + \frac{32}{\varepsilon^2} \ln\left(\frac{1}{\delta}\right)$$

$$* = \frac{32d}{\varepsilon^2} \ln(n) + \frac{32}{\varepsilon^2} \ln\left(\frac{1}{\delta}\right) + \frac{32d}{\varepsilon^2} \ln(n)$$

now we can find a numeric bound
on n

$$n \geq \frac{256}{\varepsilon^2} \ln\left(\frac{128d}{\varepsilon^2}\right) + \frac{256}{\varepsilon^2} \ln\left(\frac{1}{\delta}\right) \geq \frac{256}{\varepsilon^2} \ln(128d)$$

$$\left| \begin{array}{l} \delta < 1 \\ d \geq 1 \end{array} \right\} \geq 256 \ln(128) \geq \frac{20 \cdot 4}{d}$$

$$\Rightarrow n \geq \frac{32d}{\varepsilon^2} \ln(n) + \frac{32}{\varepsilon^2} \ln\left(\frac{1}{\delta}\right) + \frac{32d}{\varepsilon^2} (\ln(2e) + \ln(n) - \ln(d))$$

$$\geq \frac{32}{\varepsilon^2} (d \ln(n) + \ln\left(\frac{1}{\delta}\right) + d \ln(2e) + \ln(n - d \ln(d)))$$

$$\Rightarrow \frac{n\epsilon^2}{32} \geq \ln\left(\frac{2en}{d}\right)^d + \ln h + \ln\frac{1}{\delta}$$

$$e^{n\epsilon^2/32} \geq h \left(\frac{2en}{d}\right)^d \frac{1}{\delta}$$

$$\delta \geq h \left(\frac{2en}{d}\right)^d e^{-n\epsilon^2/32} //$$

Also we know from PAC learnable
for the Agnostic case ($\text{VC dim} = d < \infty$)

$$P(e_p(\text{ERM}(S)) - \min_{h \in H} e_p(h) > \epsilon) \leq h \left(\frac{2en}{d}\right)^d e^{-n\epsilon^2/32}$$

$$\Rightarrow P(e_p(\text{ERM}(S)) - \min_{h \in H} e_p(h) > \epsilon) \leq \delta //$$

Programming Assignment

- (a) See stdout when running script
(b) we can consider the classifier image of perceptron as an heat map in which pixels with high and positive values are strong characterizers of "8" and not of "0" these pixels are those of the "cross" in the middle of "8" these pixels

and pixels that are significant in the "0" digit but not in "8" will have negative values

~~look~~
those pixels

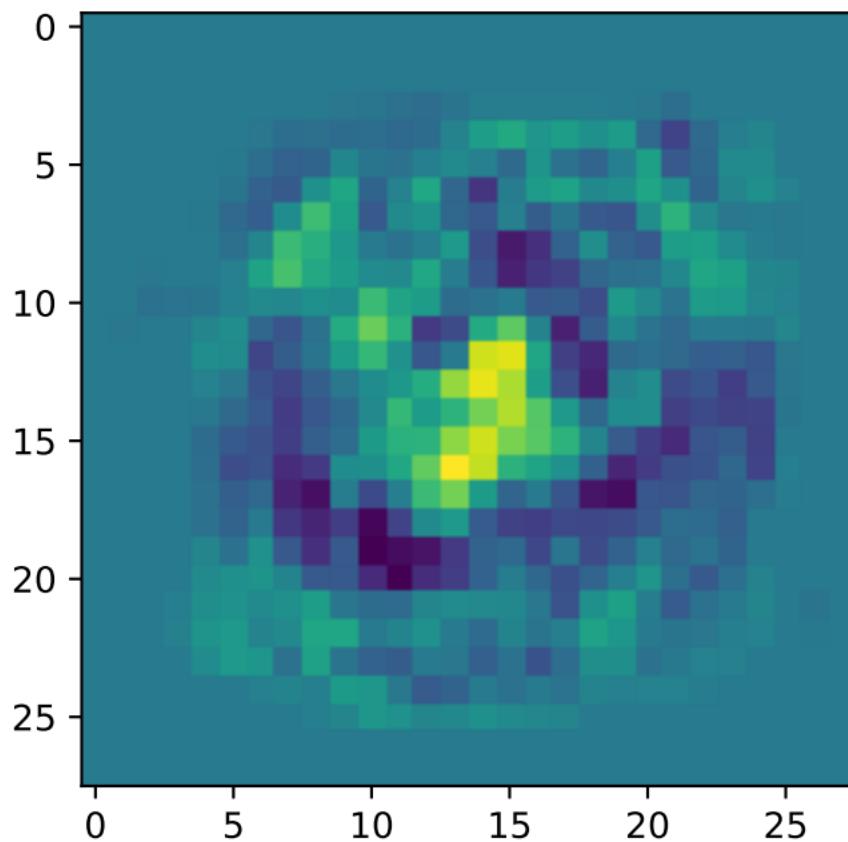
in that way, a scalar product of the classifier with "8" pictures tend to yield a positive value and with "0" pictures - a negative value - as expected from the classifier

- (c) Printed as stdout. see script or better, run it yourself

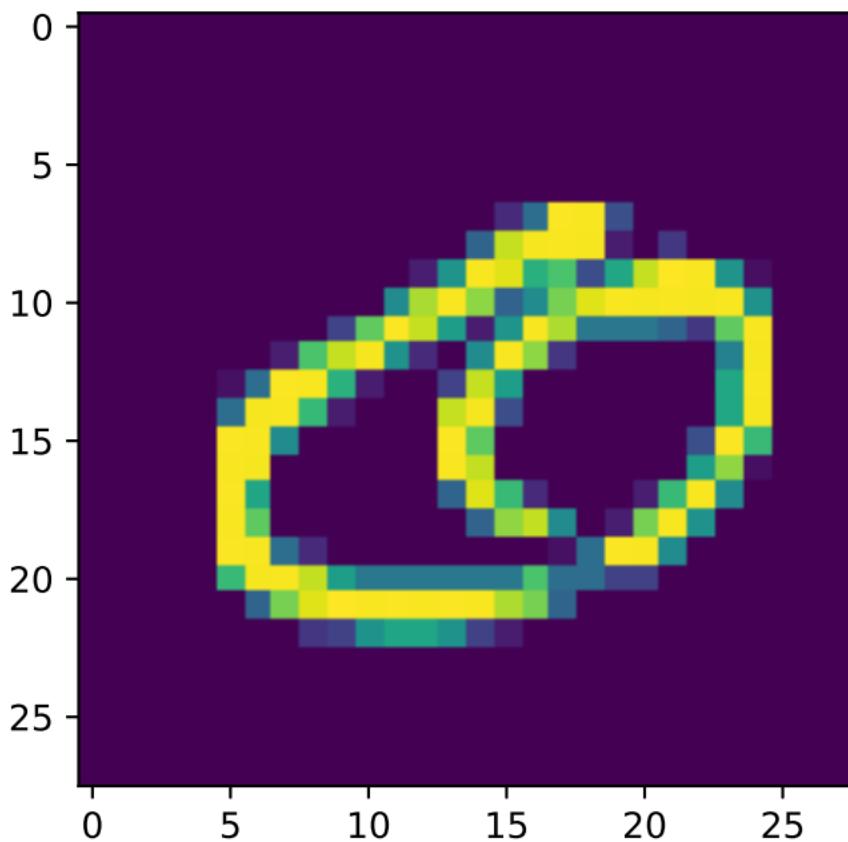
- (d) The two misclassified pictures are "0" pictures that were classified as "8". not surprising, we can see that in both there is a smudge towards the middle of the picture, why the classifier has high positive values for classifying positively

② see all results in stdout
and attached plots.

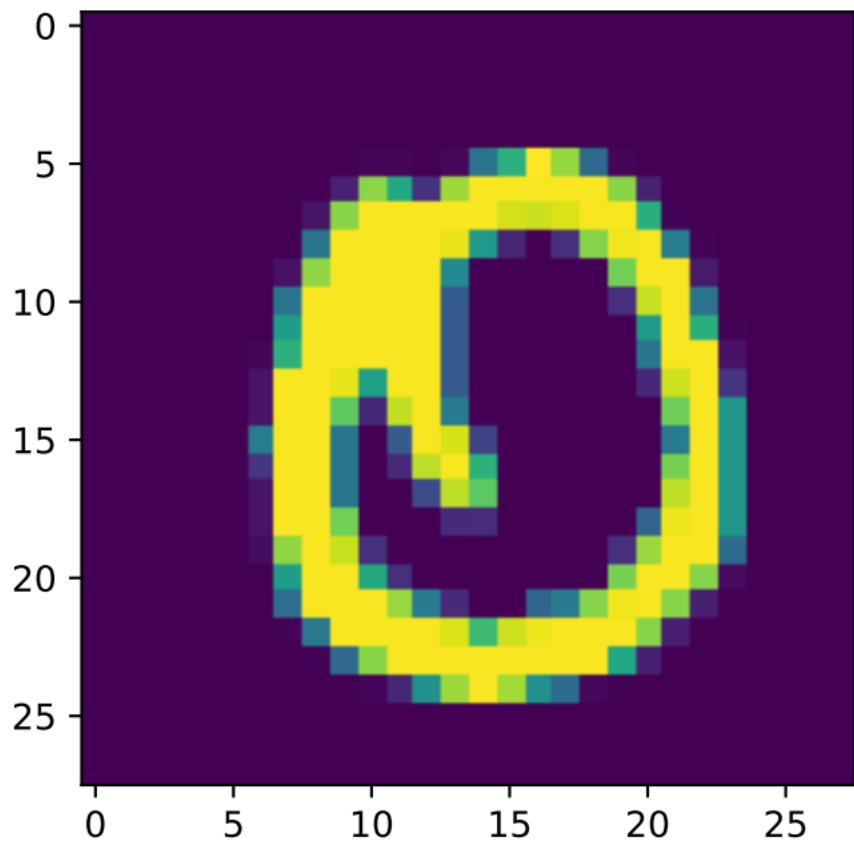
Perceptron Classifier Image from Perceptron
on Full Train Data Set

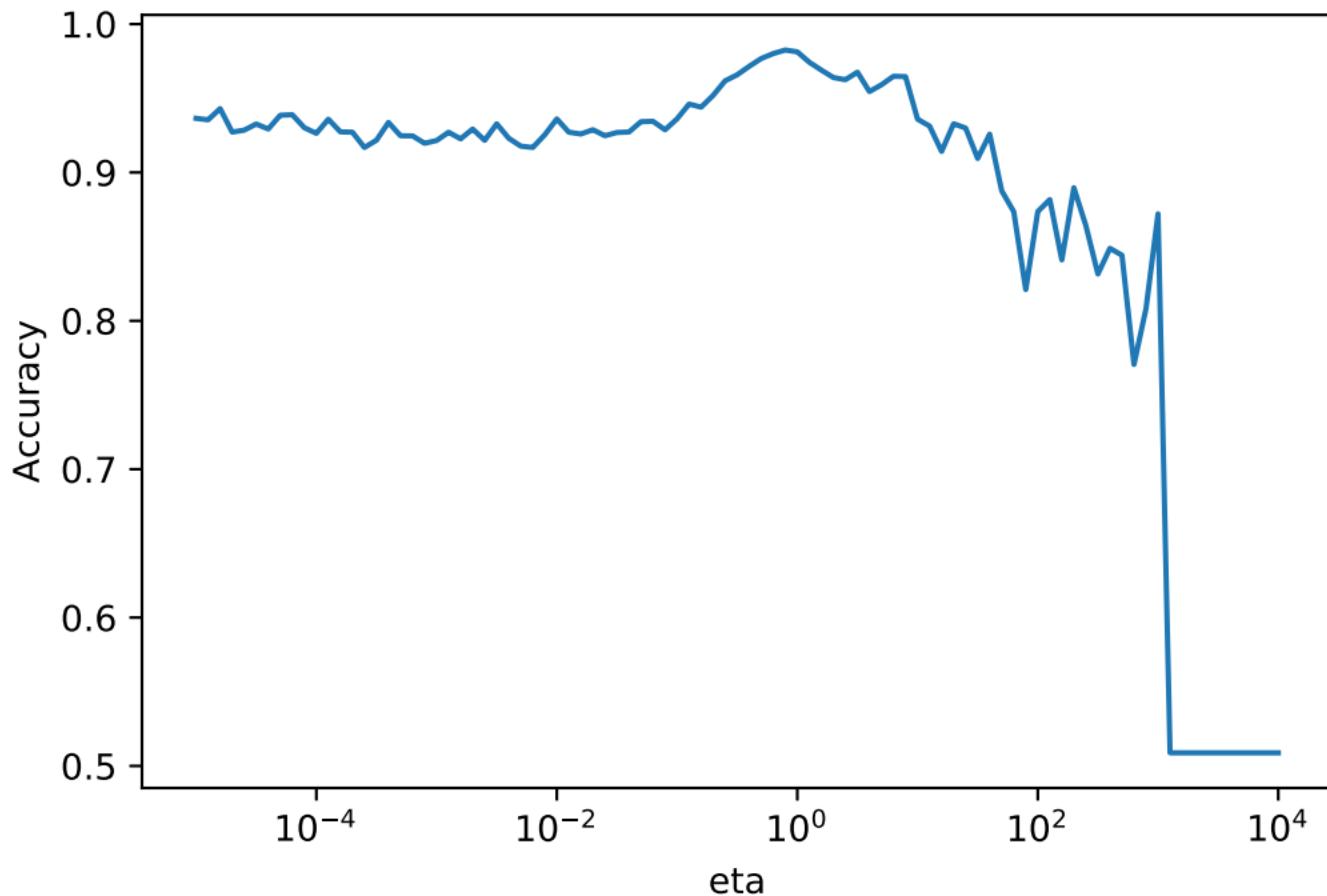


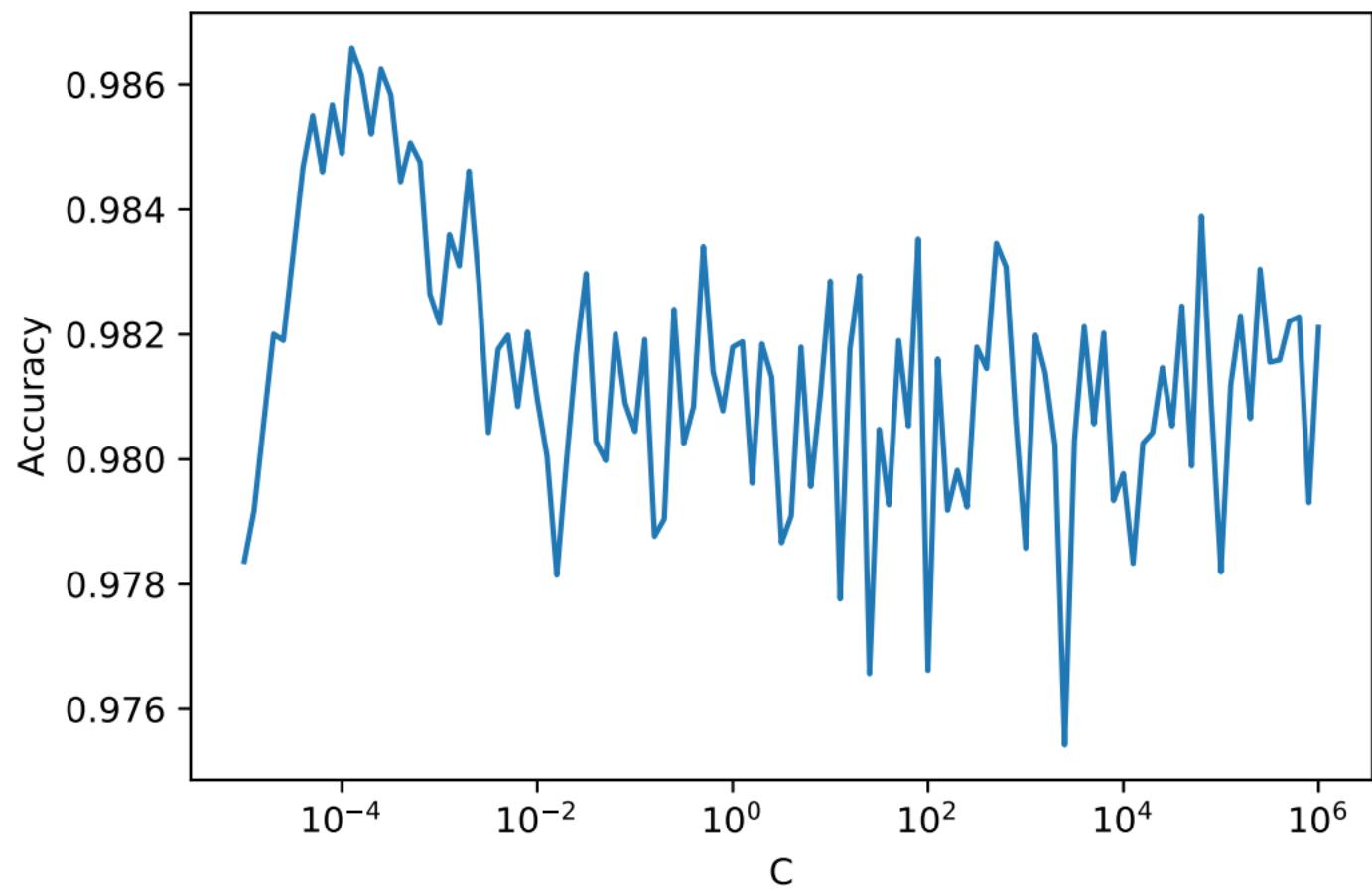
A misclassified test Image by Perceptron
on Full Train Data Set



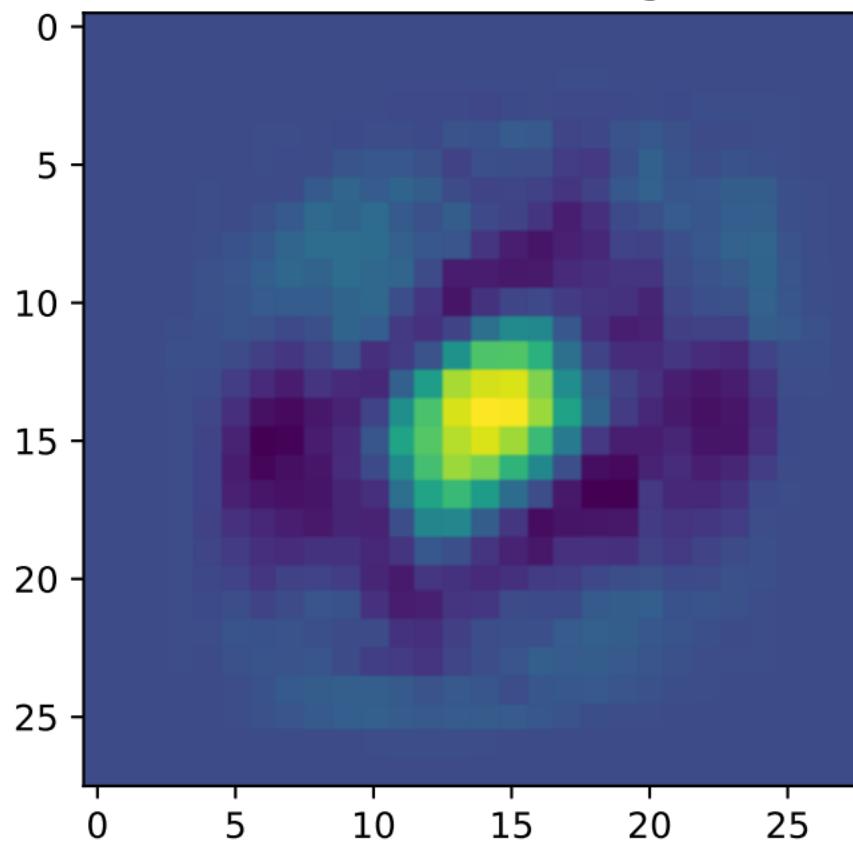
A misclassified test Image by Perceptron
on Full Train Data Set







SGD classifier Image



```
C:\Users\shosh\Google_Drive\Studies\Intro_to_ML\ex\ex3>python perceptron.py

Running perceptron for train data sizes = [ 5 10 50 100 500 1000 5000], 100 times per train data size

Accuracy results on test data:
N      Accuracy      5 percentile    95 percentile
5      0.9028        0.8332          0.9335
10     0.9096        0.8677          0.9483
50     0.9591        0.9295          0.9791
100    0.9684        0.9471          0.9826
500    0.9812        0.9611          0.9898
1000   0.9859        0.9769          0.9913
5000   0.9891        0.9805          0.9933

Running perceptron for the full train data set
Accuracy of the result classifier on the test data set is 0.9877
```

```
C:\Users\shosh\Google_Drive\Studies\Intro_to_ML\ex\ex3>python sgd.py

Running SGD for different values of eta_0, T = 1000, C = 1, iterations per eta = 10

sgd.py:56: RuntimeWarning: overflow encountered in multiply
  w = (1-eta_t)*w

Best eta is 0.794

Running SGD for different values of C, T = 1000, eta = 0.794, iterations per C = 10

Best C is 0.00025

Running SGD for C = 0.00025, T = 20000, eta = 0.794

Accuracy of best classifier is 0.993
```