	1 1
	1. GRADIENT DESCENT $f_{q(2)} = qx^{2} \text{for a 70}$
(i) we one given takes = of
	11 2 step sizes
4.	let. a=2, we consider a-step sizes
	- N = 2, N = ? Starting gradient descent at
	0.824
	$z_0 = 1$?
	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
	Soln.
	that we
	Since we arready know the value of 9=2, thon we
	proceed as below
	3
	fay = 2x2 - Tho is the thinching we going
1	be using to determined the optimal Whiting
71547	with goen step sizes
	1 - D
	=> We first find the derivative at the 9614
	Tunchin.
	=) $CH_{PU} = 4 \times -(i)$
t-	de
Ta 1) Then at Zo=1 Egn (i) will veril to
	$f(\alpha) = df(\alpha) - 4\alpha$
	Δa $\alpha = 1$
	= 1/0 = d/0 = 4/0 = 4
	da - +(0 = +
1	
-	one observation to note is that the directive
	is positive i.e = 4 > So with the we already know that the
	Optimal value solution is Jething bigger larger thus we want to go back words
	Coptinal Vaine Solution 12 Jethry bigger larger
	p. 1 9/4/24 / (4/4/24)

1 hat a
=) In this respect we use the below fuctor, that
the state of the s
$\chi_{i+1} = \chi_i - \prod_{i=1}^{n} \chi_i$
so to them of had of book on cl (=
given the 2- Her sizes of tollow
21/2016 (4)
2'ity = 2; + 21 R)
$1 \times 1 \times$
grades I a find the middle of deduction
>> Zith = 1 + 19 (or) put f (or) = 4
to be made for any formation and costs (a)
= 1 = 1 - \((4) = -1
=) we use -1 as the now a value to tead trums
Little ean (iii) as below
is the same of the same of the
=7. Ebert, me combute put gainsput raine & for at
The new & value 1:-1 , 111 personal tout
2.7
=) - (=0) = 42 = 4(-1) = -4
=) The the resulting chanceron of the 2 new
Unine will be que as
$e_{i+1} = -1 - a(-4) = 7$
=) S. Starting gradient devent will result to
a sudden increase and he diversion at family of
functions faces which is contributed by a
resulting increase in & value starting xo=1

Detming Set all step iron no for which the gradient descart will tend to Caveryo be the minimum. to top to hourston at boost of whom , or it all step trees, we need to consider the following to First we now to find the value of 1 at which the gradient desent will converge (ii) And then using that 1-vaine we determine he to value to contram that it is actually Starting at X0=0 (11) Then to determine the Set of all App Grass
we have to take opposite interval Starting but not including the 1 1-value i.e I Since we are quaranted it a positive devinating from las = 212, we know he correct to use is great by as my and when 21, for = 21 + 17 1/00 = 1 Howard Now (et of 2500 A) from (i) and alsumny 1 is notingen the though will andrew forestaling to the forest hypothesis

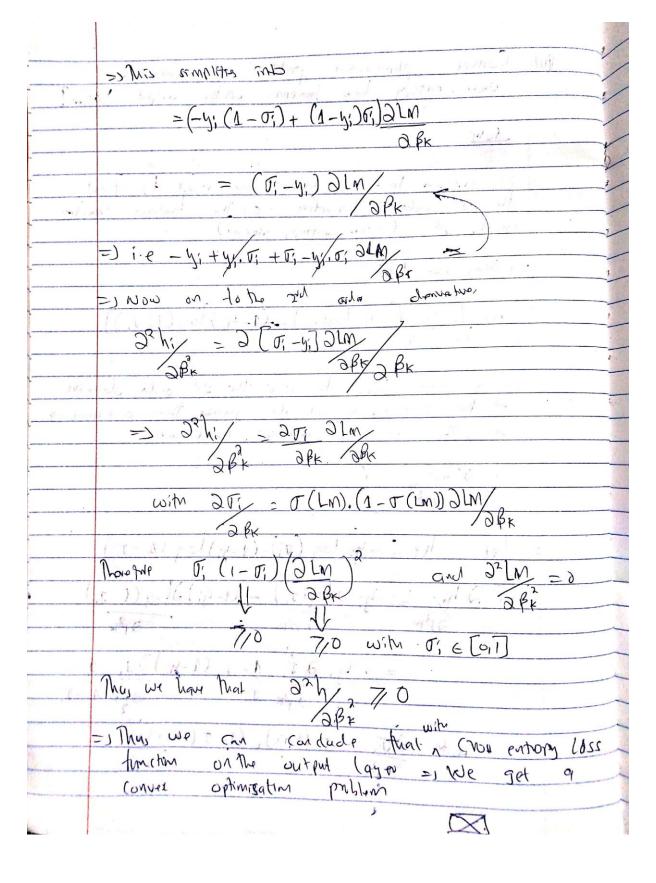
	=> 2(+1=1-1(4) -1-47000 19) (m)
-	s Now we are this value as an new x-vaine
	and substanting it in
4,-	for = 2 x2 we have
	de ou to march at any of the second
	f.(1-47) = 2 (1-41) = 2 [1-82+162]
11.6-	$2 - 160 + 320^2 = 90$
1"	and the second s
e # ($= 320^{\circ} - 160 + 2$
	-125
5	so to get Iminimum, we cet do =0
	and the state of t
	=> dgn = 6+2-16=0
	d(n)
	=> 642=16 => 1= 0.28
	where the transfer of the second
	dets détermine 2/10
	A 6.1
· /	= 1 - 4n = 0 - 1 = 0
/ <u> </u>	=> 20 + Out and all the
	a support of the land of the l
	=) So the step-1720 of (noludes, all values
1 2	other than 1=0.25 and 170
	3 3 1° 1 3: 1 sayle congres druck
	=) open-interval tum 0-27 to too
	The state of the s
	Set ut stall 1= oner our < 2 4 wy
	ungrowth is a surround a (15)
The same and	OR is Real No. System

(ii) Let 276, other how many leps, it do we
(ii) Let 276, orther now way
have 2', -2* LE
Che con the second of the seco
Chy and the number at
In other words we are to prove the number at
(farations (1) need to achieve convenience of
The gradual odescent functions
=> So suppose the trunction grow faces: R" > R
1 Harrish Alleha
gradient is Lipschitz Continue with Courtent
for any x it we man gradient descont
1 (1 (1)) (1) (1) (1) (1) (1) (1)
we yield a solution of the finalisch Station
The below tunction
3 (a) (b) (c) * 12
$f(x^{(i)}) - f(x^*) \leq x^{(i)} - x^* _2$
ation and a time and a series of the series
The take of Convergence of mis function is the
0 (1/1)
=> so with this we can easily. determine the
i iterations neaded for Conveyance
201 1/10 12 politica) J. 10 1212 1942 1/1 05 (5)
=> We consider the bound of world (201) - f(20) 5 &
Myry 1 Myles 1/107 3: -34 7 8
=> Thus with Import land of Convergence rate show
above it can only be achieved using
171 11 11 11 11
(log(1/2)) number of iterations
Comple with the site has been seen to the seen t

(1)	Salar de la company de la comp
Apr . ·	So in this since the multiclass logistic regression
	model is more powerful and able to
	represent more complex relationing between
y.	Empate and output
711. 2	in a training the war to make the same
A	=> One important key takeaway when determing
32 5 1	efficiency of a model is that we need
gdt !	to Consider the number of hidden units
	Which shall in most cases being than
By Ju	the twice the gize of the input layor
Vary of	At is want of with open day and silvert
10	=) In our case , 25 hidden units and the
100 %	number of dayer in the input layer being
	only to, this may lead to the model being
	overtitled and with the tack that no
· Jan	regularization on that model => simply
of	means that this model will work only
- Jan	ettracenty on the training data but will
sp 1.	
	output of the actual test data
fre y	
904	
j. v	kept between the cize of the input land
- Fly	with the same of t
	12 not the case with the 25-hidden units madel
	thus making it more prone to overtiting
	and hence less power tull than the
	malti das alogitic regionim
	The multi clau logitic regression can done with
	explitionly number at classes hence out enting
	Values for each class with no overfitting

212 A James of Harrie of => meaning that the model works better on both training and testing data with high according wither 4 hodden units, the model will be a bit powerful than the multiclass legister regression in the sonse that, the number of hedder units is kept in between the size of the => Also during training this model will out-partium Multidas Cogritic regiension because at the loner Number of hidden units used in the hidden layer making it. pretty tast. and powerful =) we also need to mention that the choosen hidden Lay units makes the model not being undertitled => implying that to is going to work efficiently on both. The training and terms data" with high according and fout speal =) Thus model with 4 hidden units is better in this case than multiglass logistic regression The fact lite takes lesser time to train model and will obtains high accuracy with The girtual datagen two wast could the

(ii)	Convex optimisation problem if we use
	crow entury loss tweeter on the output layor?
	(a) ((1) ((a - b)) + ((0 - b)) ((a -
	Soly
= 3	To answer the problem we need to trud.
	The 2th order derivetive of the Salow function
	with B (crop entury function)
	= 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
	= 1 Mas we have.
	-dank.
	h= [5 [-4, log(vi) (1-4,) log(1-vi)]
	N 120 0 1111-111 (- 1111) (- 1111)
	13 2942
	=> 20 ms need to twy to zzy arga garacter
	of the trucky, so it its greated than or equal to
	Zoro, then we can make conduing
	010 . 116
	23/ 2= 7
	14/2Por 15/10/10 - 1 (4) 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
	=> let h; = -4; log (T;)-(1-4;) log (1-T;)
4	2/2/5/ 1/2/2/5/
	=) 2 hi/ = -y; 2 log(vi) - (1-yi) 2 log (1-vi)
	Japk Japk
	$= -4', \partial \tau; \Delta + (1-4!) \partial \tau; \Delta$
	$= -y', \partial \tau; 1 + (1-y;) \partial \tau; 1$
	2 PK 1-1;
ha & l	4) M 30; = T (LM), (1-T (LM) LM
11.67	
No. of Concession, Name of Street, or other party of the Concession, Name of Street, or other pa	
The second second second plans	- DC)



	3. Digit dassitication using MLPs
(1)	independent the rise of the network and
	Willes Officiary
1000	-3 m/10/100
AND 112	(1) Umy binary Encally instead of one-hol-Encally
	13 reduce Norther of narrow in the middle land
72. 17	khat to think about the uggethe?
()	July allege the world
	co to answer the april are drougher one down
	implement the exact never notions too as
	described above and Compaged it alw
بسيا	dosabod above with binary encolong lover
	Neurons in the middle layer and compared
	ils patomanco with a NN with une-hot-
J	Pencoding israeds surveyed it will be
,	Harden = >00 als misseles i labore will present
ward	So now based on the experiment Conducted
	on the implemented moder with binary encoding
	and level Newm => this does not actually
	reduce the size and efficiency in neare
	= 1 By evaluating the two models pertomence 1) model using binary encoding had a accuracy of
1337	(i) model using binary encoding had a accuracy of
	900 079
	(ii) Model using one-hot encoling had an
	accuracy of 96%
	=) From these results it is clear that these
	suggestions actually do not result into a
	More efficient model
,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,	Thus a model with one-hot-encoding is much more
	improved and high officient.

Question 3B) See Code attached and only a from a => So from the result obtained , the run madel a courage on the model we obtained a 96% whole and 92% accuracy using one-hot-encoding white and and on he model usmay an addittonal layer when is relatively dose Nertication See attached code for) -2. happens when the NN above is trained . What. without adding the last layor => From the experimental observations is that we are actually training this model directly with a lower getting a lesser ettricent model arrivary of 87% See attached code