

CHALMERS

EXAMINATION / TENTAMEN

Course code/kurskod	Course name/kursnamn		
Anonymous code Anonym kod	Examination date Tentamensdatum	Number of pages Antal blad	Grade Betyg
Dit 821 369	Software engineering for AI 2022-10-26	11	VG

* I confirm that I've no mobile or other similar electronic equipment available during the examination.
 Jag intygar att jag inte har mobiltelefon eller annan liknande elektronisk utrustning tillgänglig under
 eximinationen.

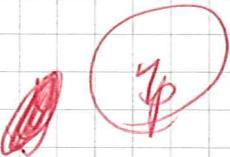
Solved task Behandlade uppgifter	Points per task Poäng på uppgiften	Observe: Areas with bold contour are to completed by the teacher. Anmärkning: Rutor inom bred kontur ifylls av lärare.
No/nr		
1	X 4.5	
2	X 4	
3	X 5	
4	X 3,5 4	
5		
6		
7		
8		
9		
10		
11		
12		
13		
14		
15		
16		
17		
Bonus poäng	1	
Total examination points Summa poäng	12,5	11,5

CHALMERS DIT 821	Anonymous code Anonym kod 369	Points for question (to be filled in by teacher) Poäng på uppgiften (fylls av lärare)	Consecutive page no. Löpande sid nr 4
			Question no. Uppgift nr 2

a) The decision boundary equation is:-

$$Y = g(6 + (-5X) + 1 \times X^2)$$

$$Y = g(6 - 5X + X^2)$$



b) to solve this question we will choose $x^{n \times m}$ that output the maximum probability for y_1, y_2, y_3 the labels for the three new labels is

$$(x_1, y_1), (x_2, y_3), (x_3, y_1)$$

$$= (x_1, 0), (x_2, 2), (x_3, 1)$$



c) The certain parameters is the "K" number of clusters it group and classify the data to this no specified number of classes or cluster.



~~k) k-means algorithm continuation~~

~~Repetit~~

The higher K number means we add more clusters.

d) Continuation of k-means algorithm :-

Repeat {

for $i=1$ to m

$C^i = \text{average of points closest to } x^i$
Compute distance and assign the
closest x^i .

for $k=1$ to K

$\bar{x}_k = \text{assign average(means) of points of}$
 $\text{all } x^i \text{ closest to this Centroid.}$

{

7p

c)

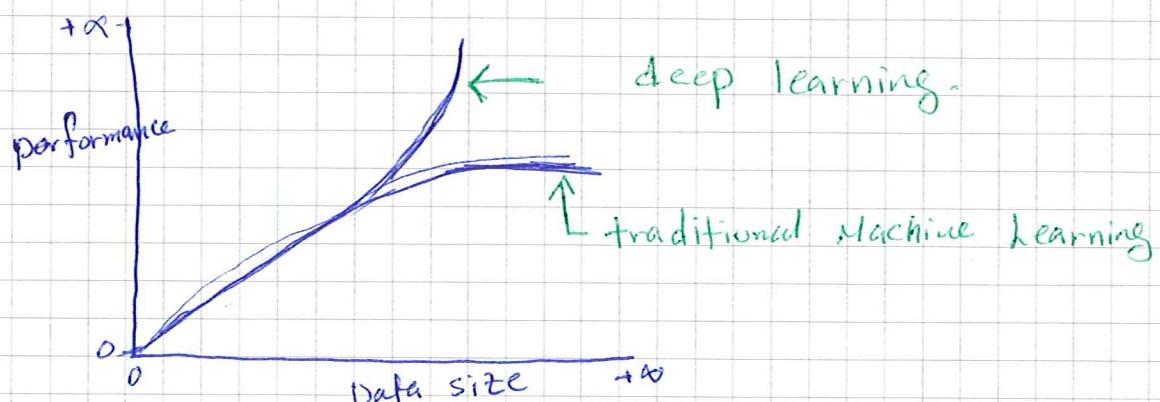
Sum of the distances between points and
their centroids squared divided by the sample
size or the number of training set.

The goal is to find The minimum distances
between cluster centroids and points which
result in better clustering -

1p

a) In traditional ML we must perform feature engineering and feature extraction while in Neural Network is not necessary. Can done automatically and feature selection through the NN

b) Look at the graph below. In traditional ML perform good when the data size is small but then it will have performance parallel to x-axis in Deep learning. The performance increase as the data increase.

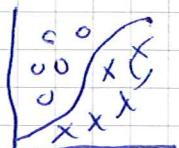
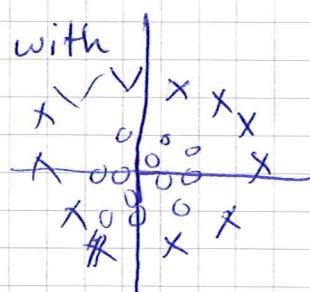


c) Deep learning deals with very complex problems that traditional can't as in traditional a Logistic Regression deal with one decision boundary but it can not deal with

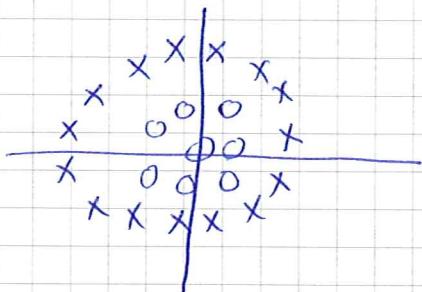
this problems →

where the decision

boundary is non linear.



b) non-linearity allow neural network to deal with very complex data "see the graph" such as:-



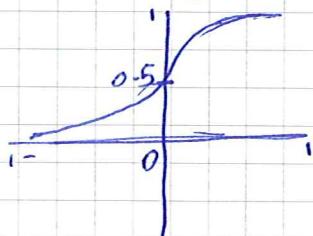
it uses activation function to tackle this problem

where there is no one decision boundary.

* function that is used is Sigmoid Activation function

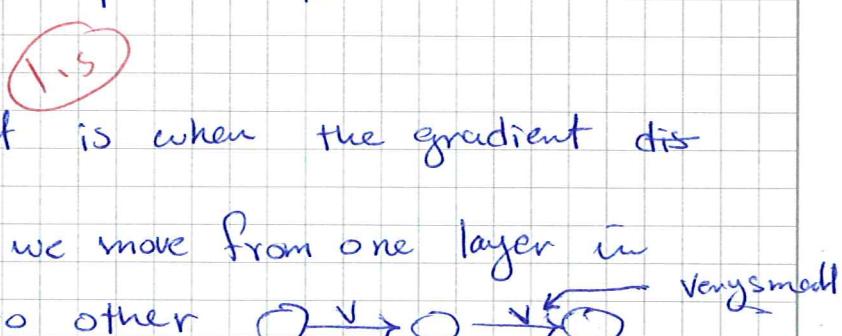
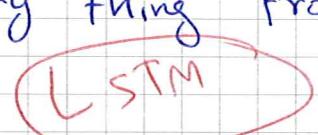
$$g(z) = \frac{1}{1+e^{-z}}$$

$$z = h(x)$$



c) it is smaller because we want to detect only the important feature and to increase decrease ^{0.5} model complexity.

and to detect certain features that exists ~~in~~ in certain location in the image.

CHALMERS	Anonymous code	Points for question (to be filled in by teacher)	Consecutive page no. Löpande sid nr
	Anonym kod <i>Dit 821 369</i>	Poäng på uppgiften (ifyller av lärare)	Question no. Uppgift nr <i>3</i>
		(1)	
	d) i) deal with spatial variance to reduce distortion in image.		
	2) reduce model complexity by reducing number of parameters which increase model performance.		
?:	3) Capture the most important feature and maximize it to better detection in the model and eliminate non important informations.		
	e) Vanishing gradient is when the gradient disappears when we move from one layer in Neural Network to other  Very small		
	RNN solve it by using sigmoid gate and point wise multiplication. when sigmoid is equal to 1 it means allow all the information to pass to next and when is 0 it means block every thing from flowing to next layer.		
	 LSTM		

a)

b) wrong measures of the data could be that can appear during data collection. and systems breaks as well-

2) Social biased data which lead to biased model to for example gender or race.

1

b) it is important to Categorical data to numerical data for the model to be applicable to process. ^{which one?}

~~one hot~~ one example is the house data

Size	Pool Area	Parking area	Price
10	0	0	1000
5	5	5	2000
8	0	0	5000
9	8	8	10000



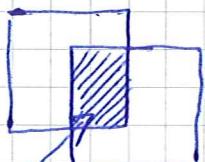
Size	Pool Area	Parking Area	Price
NA	NA	NA	
medium	NA	NA	
NA	NA	NA	
large	large	large	

One-hot encoding is used to map categorical data into numeric data either by assigning binary values 0 to not available or no and 1 to available or yes or by specifying range for example if $x = 0$ if pool area is NA and $x = 5$ if pool area is medium. ✓ ~~5~~
still 1

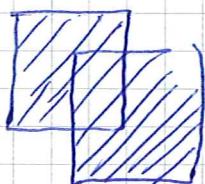
c) the formula is to test Inter-Agreement Annotator

between two or more annotators simply by calculating the Area of two intersections of two boxes and divide it by area of union of two boxes:-

* The boxes is the specific Area of the Image determined by Annotator and labeled like Car Driving way, house, etc



intersection of two boxes



Union of 2 boxes

the output will be score between 0 and 1
0 for no matching , 1 represent total Agreement

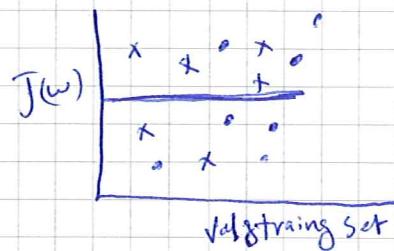
0.4 represent poor match, 0.7 = good match
and 0.9 Excellent Agreement.

1

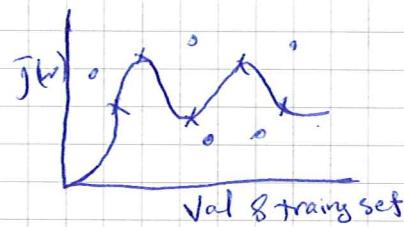
CHALMERS	Anonymous code	Points for question (to be filled in by teacher)	Consecutive page no. Löpande sid nr
	Anonym kod <i>Dit821 - 369</i>	Poäng på uppgiften (ifylls av lärare) <i>1</i>	Question no. Uppgift nr <i>4</i>

d) it's not always applicable ~~as if~~ in most cases we could generate metrics from the requirement for example the house prediction the price should fall between minimum value and maximum value for price but in some cases we can not specify metrics & for evaluation from the requirement like algorithm that predict human behaviour also to measure measure how bias the model is- ✓

- a) Under fitting is when the model is too simple if it can't generalize well in the training data and validation set.
- b) Over fitting is when the model is too complex and it fits perfectly the training data but it produce high error in validation set.

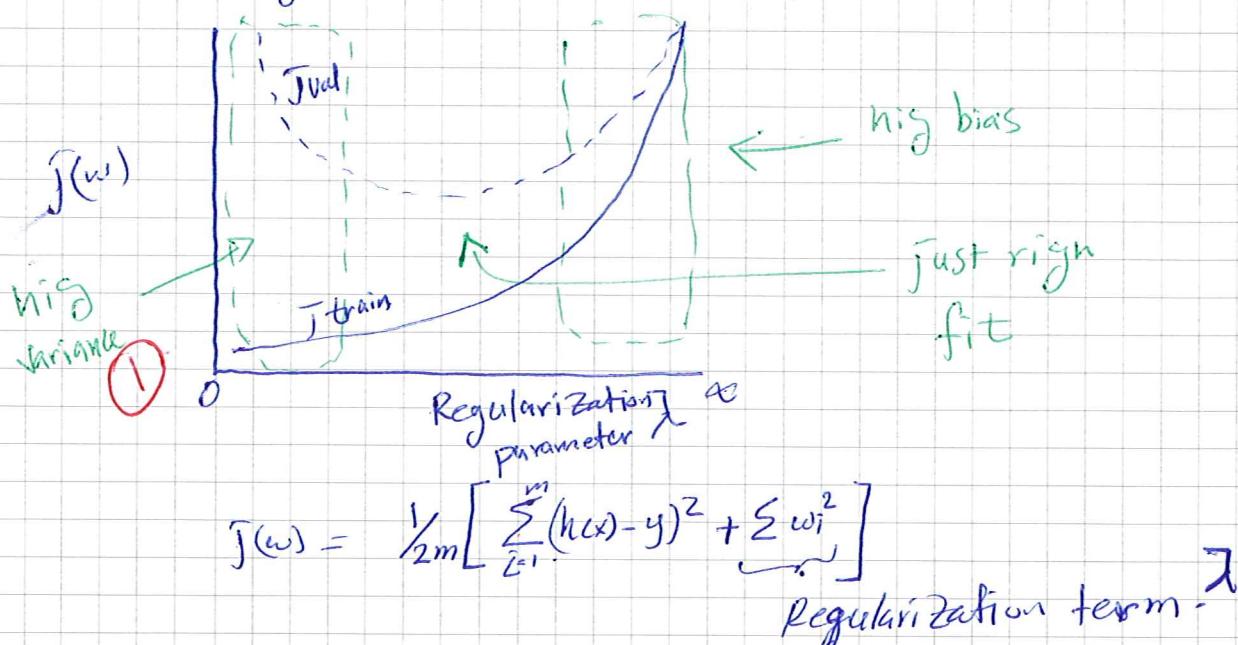


(1)



x training set
o validation set

- b) Regularization help avoid overfitting by minimizing the effect of the high order features x^4, x^5, \dots
- So the model will not just fit training data but it will generalize well in the validation set.



it will penalize the model with small value related to model complexity. the more complex model the bigger lambda is needed to better fitting the data

c) gradient decent equation

$$\hat{J}(\omega) = \omega - \alpha/m \sum_{i=1}^m (h(x_i) - y_i) \cdot x_i$$

- Normal equation

$$\omega = (X^T X)^{-1} X^T y$$

- in gradient decent we must specify number of iterations while in normal equation not required.
- in gradient decent we must choose alpha while in normal equation not required
- gradient decent work well with big Data ~~data few~~
- ① while normal equation does not work with ~~not work with~~ big Data ($X^T X$) is computationally expensive and the inverse in $(X^T X)^{-1}$ is more expensive and some time is not applicable-

d) I would prefer gradient decent as it's big Data and the normal equation is not

② work well in big Data $m=50 n=200,000$

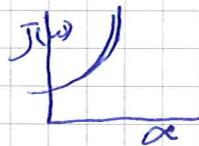
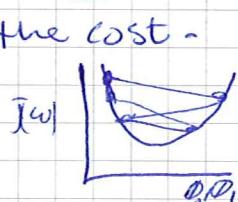
- if if we look at the equation $(X^T X)^{-1}$, the X transpose is computationally expensive and the inverse may not be applicable for ~~for~~ ^{for} ~~not~~ ^{not} big Data.

~~no. of~~ ^{is more}
~~not~~ ^{not} ~~detail~~ ^{detail}

e) a) false if we specify α is large value

the gradient decent may not converge and we might over shoot the optimal points that result

in increasing of the cost -



α is large

b) False the

the reason is that each one has different effect of the cost function w_0 may decrease faster or slower than w_1 , since we might have local optimal.