

Note:

- During the attendance check a sticker containing a unique code will be put on this exam.
- This code contains a unique number that associates this exam with your registration number.
- This number is printed both next to the code and to the signature field in the attendance check list.

Introduction to Deep Learning

Exam: IN2346 / Endterm

Examiner: Prof. Leal-Taixé and Prof. Nießner

Date: Tuesday 11th August, 2020

Time: 08:00 – 09:30

P 1 P 2 P 3 P 4 P 5 P 6 P 7 P 8

--	--	--	--	--	--	--	--

Left room from _____ to _____

from _____ to _____

Early submission at _____

Notes _____





X

X

X

X

X

X

X

IN-I2DL-1-20200811-E0134-02

IN-I2DL-1-20200811-E0134-02

IN-I2DL-1-20200811-E0134-02





Endterm

Introduction to Deep Learning

Prof. Leal-Taixé and Prof. Nießner
Chair of Visual Computing & Artificial Intelligence
Department of Informatics
Technical University of Munich

Tuesday 11th August, 2020
08:00 – 09:30

Working instructions

- This exam consists of **20 pages** with a total of **8 problems**.
Please make sure now that you received a complete copy of the exam.
- The total amount of achievable credits in this exam is 90 credits.
- Detaching pages from the exam is prohibited.
- Allowed resources: **none**
- Do not write with red or green colors nor use pencils.
- Physically turn off all electronic devices, put them into your bag and close the bag.
- If you need additional space for a question, use the additional pages in the back and properly note that you are using additional space in the question's solution box.



Exam empty





Problem 1 Multiple Choice Questions: (18 credits)

- For all multiple choice questions any number of answers, i.e. either zero (!), one, all or multiple answers can be correct.
- For each question, you'll receive 2 points if all boxes are answered correctly (i.e. correct answers are checked, wrong answers are not checked) and 0 otherwise.

How to Check a Box:

- Please **cross** the respective box: (interpreted as **checked**)
- If you change your mind, please **fill** the box: (interpreted as **not checked**)
- If you change your mind again, please place a cross to the left side of the box: (interpreted as **checked**)

a) Which of the following statements regarding successful ImageNet-classification architectures are correct?

- ResNet18 has more parameters than VGG16.
- VGG16 only uses convolutional layers.
- InceptionV3 uses filters of different kernel sizes.
- AlexNet uses filters of different kernel sizes.
(Handwritten note: A red circle is drawn around the word "different" in the sentence above.)

b) You train a neural network and the loss diverges. What are reasonable things to do? (check all that apply)

- Decrease the learning rate.
- Add dropout.
- Increase the number of parameters.
- Try a different optimizer.

c) What is the correct order of operations for an optimization with gradient descent?

- (a) Update the network weights to minimize the loss.
- (b) Calculate the difference between the predicted and target value.
- (c) Iteratively repeat the procedure until convergence.
- (d) Compute a forward pass.
- (e) Initialize the neural network weights.

- eadbc
- ebadc
- edbac
- bcdea

e d b a c

d) Consider a simple convolutional neural network with a single convolutional layer. Which of the following statements is true about this network?

- It is translation invariant.
- All input nodes are connected to all output nodes.
- It is scale-invariant.
- It is rotation invariant.



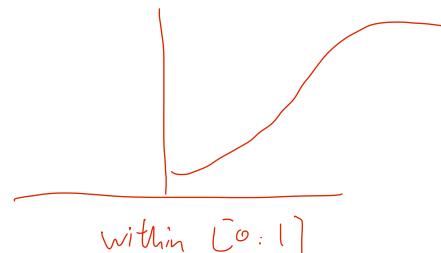


e) Which of the following activation functions can lead to vanishing gradients?

- ReLU.
- Leaky Relu.
- Tanh.
- Sigmoid.

f) Logistic regression (check all that apply).

- Allows to perform binary classification.
- Has a discrete output space.
- Uses cross-entropy loss.
- Can be seen as a 1-layer neural network.



g) A sigmoid layer

- has a learnable parameter.
- maps surjectively to values in $(-1, 1)$, i.e., hits all values in that interval.
- is continuous and differentiable everywhere.
- cannot be used during backpropagation.

h) Your training error does not decrease. What could be wrong?

- Dropout probability not high enough.
- Bad initialization.
- Too much regularization.
- Learning rate is too high.

i) Which of the following have trainable parameters? (check all that apply)

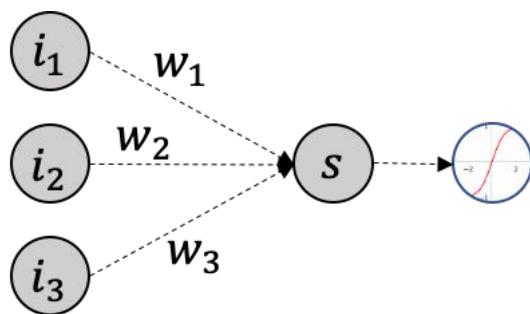
- Max pooling
- Dropout
- Batch normalization
- Leaky ReLU





Problem 2 Activation Functions and Weight Initialization (8 credits)

For your first job, you have to set up a neural network but you have some issue with its weight initialization. You remember from your I2DL lecture that you can sample the weights from a zero-centered normal distribution, but you can't remember which variance to use. Therefore, you set up a small network and try some numbers. You initialize the weights one time with $\text{Var}(\mathbf{w}) = 0.02$ and one time with $\text{Var}(\mathbf{w}) = 1.0$:



Inputs:

- $i_1 = 2, i_2 = -4, i_3 = 1$

$\text{Var}(\mathbf{w}) = 0.02$:

- $w_1 = 0.05, w_2 = 0.025, w_3 = -0.03$

$\text{Var}(\mathbf{w}) = 1.0$:

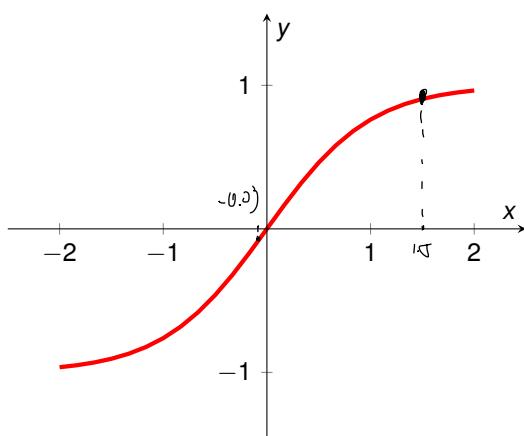
- $w_1 = 1.0, w_2 = 0.5, w_3 = 1.5$

$$\begin{aligned}
 S_1 &= w_1 \cdot i_1 + w_2 \cdot i_2 + w_3 \cdot i_3 \\
 &= 0.05 \cdot 2 + 0.025 \cdot (-4) + (-0.03) \cdot 1 \\
 &= 0.1 - 0.1 - 0.03 \\
 &= -0.03
 \end{aligned}$$

$$\begin{aligned}
 S_2 &= w_1 \cdot i_1 + w_2 \cdot i_2 + w_3 \cdot i_3 \\
 &= 1.0 \cdot 2 + 0.5 \cdot (-4) + 1.5 \cdot 1 \\
 &= 2 - 2 + 1.5 = 1.5
 \end{aligned}$$

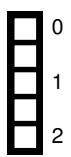
a) Compute a forward pass for each set of weights and draw the results of the linear layer in the Figure of the tanh plot. You don't need to compute the tanh.



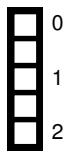


Because of the activation (tanh) layer, when input is large, gradient is equal to 0 called saturated gradient

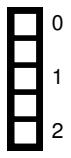
b) Using the results above, explain what problems can arise during backpropagation of deep neural networks when initializing the weights with too small and too large variance. Also, explain the root of these problems.



c) Which initialization scheme did you learn in the lecture that tackles these problems? What does this initialization try to achieve in the activations of deep layers of the neural network?



d) After switching from tanh to ReLU activation functions, one of your initial problems occurs again. Why does this happen? How can you modify the initialization scheme proposed in c) to adjust it for this new non-linearity?



C: Xavier or kaiming Initialization

They keep the gradient in the decay number not to be converges to 0

D: Because Xavier is not for ReLU, ReLU can't have the ~~input~~ output
use kaiming Initiation $\text{variance} = \frac{2}{n}$ instead of $\frac{1}{n}$



b. Small variance \Rightarrow In deep network, output is close to 0, the gradient also become very small \Rightarrow gradient vanishing gradient by variance \Rightarrow gradient saturates, gradient is also very small

c. Xavier Initialization
keep the variance of the output equal to input



Problem 3 Batch Normalization and Computation Graphs (6 credits)

For an input vector \mathbf{x} as well as variables γ and β the general formula of batch normalization is given by

b) Undo the Normalization we needed

c) BN calculate the mean and var sum
mini-batch

$$\hat{\mathbf{x}} = \frac{\mathbf{x} - E[\mathbf{x}]}{\sqrt{Var[\mathbf{x}]}}$$

a) speed up the training
prevent overfitting

Test: Only use the float mean and var from training

- 0 a) Why would one want to apply batch normalization in a neural network?

1

- 0 b) Why are γ and β needed in the batch normalization formula?

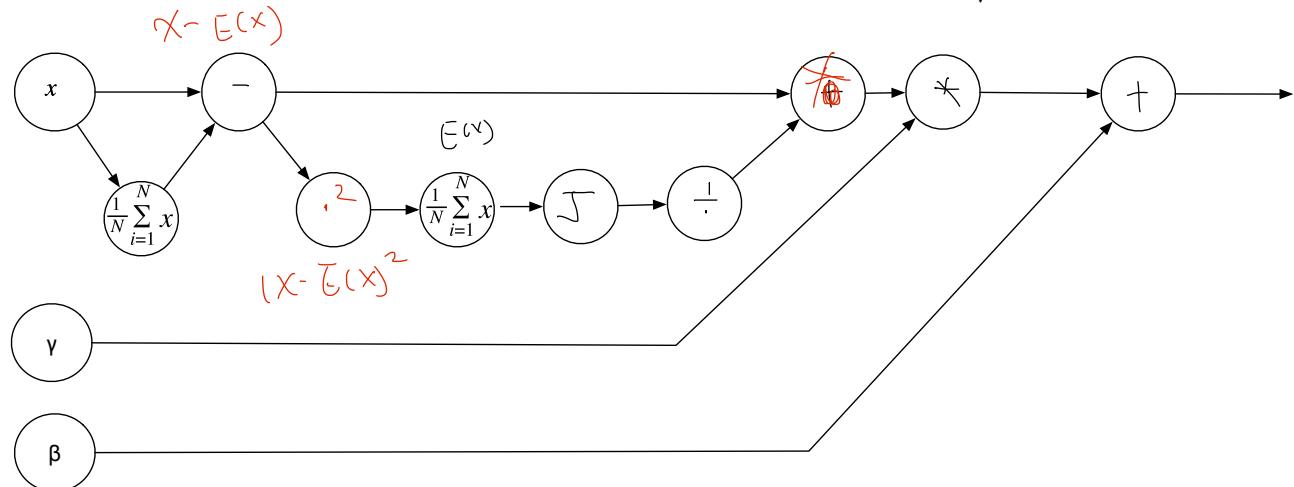
1

- 0 c) How is a batch normalization layer applied at training (1p) and at test (1p) time?

1

2

- 0 d) Computational graph of a batch normalization layer. Fill out the nodes (circles) of the following computational graph. Each node can consist of one of the following operations $+$, $-$, $*$, 2 , $\sqrt{}$, $\frac{1}{\cdot}$.



a. Reducing covariate shifted
Make Activation not die
Mimic normalization of data in each layer \rightarrow faster and more stable training

b. Allow NN to undo / shifted + scaled the normalization

c. Training : compare to mean and variance from the mini-batches
Store a weight average

Test : Use exponentially weight change mean and variance on test sample



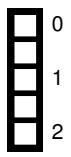
Problem 4 Convolutional Neural Networks and Receptive Field (12 credits)

A friend of yours asked for a quick review of convolutional neural networks. As he has some background in computer graphics, you start by explaining previous uses of convolutional layers.

- a) You are given a two dimensional input (e.g., a grayscale image). Consider the following convolutional kernels

$$C_1 = \frac{1}{9} \cdot \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix}, \quad C_2 = \begin{pmatrix} 1 & -1 \\ 1 & -1 \end{pmatrix}.$$

*(1) work graph dark Box mean blur
(2) detect the vertical edge and highlight the gray graph*



What are the effects of the filter kernels C_1 and C_2 when applied to the image?

After showing him some results of a trained network, he immediately wants to use them and starts building a model in Pytorch. However, he is unsure about the layer sizes so you quickly help him out.

- b) Given a Convolution Layer in a network with 5 filters, filter size of 7, a stride of 3, and a padding of 1. For an input feature map of $26 \times 26 \times 26$, what is the output dimensionality after applying the Convolution Layer to the input?

- c) You are given a convolutional layer with 4 filters, kernel size 5, stride 1, and no padding that operates on an RGB image.

1. What is the shape of its weight tensor?



2. Name all dimensions of your weight tensor.

Width

Height



Now that he knows how to combine convolutional layers, he wonders how deep his network should be. After some thinking, you illustrate the concept of receptive field to him by these two examples. For the following two questions, consider a grayscale 224×224 image as network input.

- d) A convolutional neural network consists of 3 consecutive 3×3 convolutional layers with stride 1 and no padding. How large is the receptive field of a feature in the last layer of this network?

$$\frac{F_{in} + 2P - N}{S} + 1 = F_{out}$$

$$F_{in} = F_{out} + 2$$

$$F_3 = 3$$

$$F_2 = 5$$

$$F_1 = 7$$

$\therefore 7 \times 7$





1x1 convolution layer

equal full input size

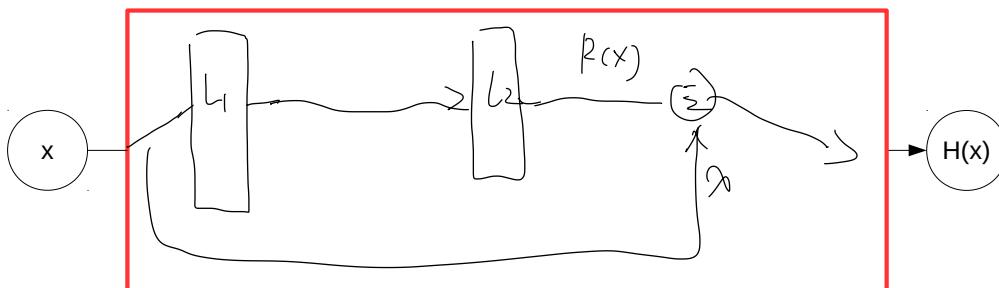
- 0 e) Consider a network consisting of a single layer.

1 1. What layer choice has a receptive field of 1?

2 2. What layer has a receptive field of the full image input?

Fully connected layer
or conv/pool layer with kernel size 1

Blindly, he stacks 10 convolutional layers together to solve his task. However, the gradients seem to vanish and he can't seem to be able to train the network. You remember from your lecture that ResNet blocks were designed for these purposes.



ResNet pass you after to the input data not the layer and the extant sum open
problem causing stop gradient flow

- 0 f) Draw a ResNet block in the image above (1p) containing two linear layers, which you can represent by l_1 and l_2 . For simplicity, you don't need to draw any non-linearities. Why does such a block improve the vanishing gradient problem in deep neural networks (1p)?

- 0 g) For your above drawing, given the partial derivative of the residual block $R(x) = l_2(l_1(x))$ as $\frac{\partial R(x)}{\partial x} = r$, calculate $\frac{\partial H(x)}{\partial x}$.

$$\frac{\partial H(x)}{\partial x} = \frac{\partial R(x)}{\partial x} + \frac{\partial x}{\partial x} = r + 1$$





Problem 5 Training a Neural Network (15 credits)

A team of architects approaches you for your deep learning expertise. They have collected nearly 5,000 hand-labeled RGB images and want to build a model to classify the buildings into their different architectural styles. Now they want to classify images of architectures into 3 classes depending on their style:



Islamic



Baroque



Soochow

- a) How and in what parts would you split your dataset to build your deep learning classifier? Formulate your answer in percentages.



- b) After visually inspecting the different splits in the dataset, you realize that the training set only contains pictures taken during the day, whereas the validation set only has pictures taken at night. Explain what is the issue and how you would correct it.



- c) As you train your model, you realize that you do not have enough data. Unfortunately, the architects are unable to collect more data so you have to temper the data. Provide 4 data augmentation techniques that can be used to overcome the shortage of data.



60
80% training set 20% test set
20% val

The division
of the dataset
is not same,
it will cause
prob prob.
we do will in
val but not
in training set

Resign the dataset mix training
shuffle split again
bad generalization
high val - error

Filipiny

Crop
Pin
Gaussian blur

Sharpen

rotation.

|@ Adding noise



- 0 d) You now start training your network using Gradient Descent (GD) on the entire training data. What might be a problem of Gradient Descent concerning saddle points or local minima of the cost function (1pt)?
1

- 0 e) While training your classifier you experience that loss only slowly converges and always plateaus independent of the used learning rate. Now you want to use Stochastic Gradient Descent (SGD) instead of Gradient Descent (GD). What is an advantage of SGD compared to GD in dealing with saddle points?
1

- 0 f) Explain the concept behind momentum in SGD
1

- 0 g) Why would one want to use larger mini-batches in SGD?
1

- 0 h) Why do we usually use small mini-batches in practice?
1

- 0 i) There exists a whole zoo of different optimizers. Name an optimizer that uses both first and second order momentum
1

- 0 j) Choosing a reasonable learning rate is not easy.

1. Name a problem that will result from using a learning rate that is too high (1p).
2. Name a problem that will arise from using a learning rate that is too low (1p)?



IN-I2DL-1-20200811-E0134-12

d. Saddle point is the point that gradient = 0.
Lak minimum
GD may have the ability to escape this point
thus the model will think it's past trained
get stuck

e. Because of the noisy update, SGD may escape the saddle point
"Stochasticity"

f. First-order Momentum
Introduce a "velocity" to the Optimizer process, which can
cancel the ~~temp~~ of ~~converge~~
Avoid getting stuck in saddle point
Stable gradient ← make gradient less noisy
if layer weight norm in any iteration, model can ~~temp~~ more easily
it will speed up the training process faster convergence

h. because practically there are many many data, use large small
batch is common very slow than batch
GPU limited

i. Adam

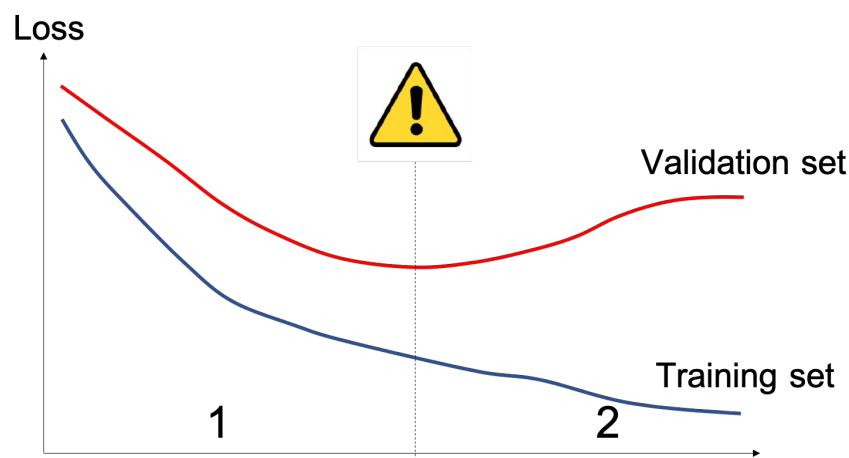
j. shoot out, miss the \oplus minima or maxima
Low density doesn't converge to optimal solution
diverge

Slow converge,
may not converge to optimal soln
or slow

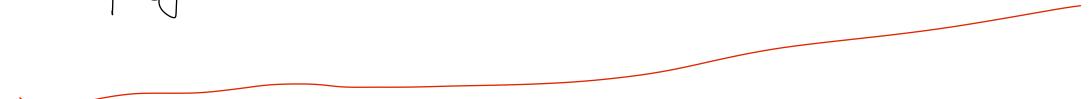


0
1
2

- k) Finally you plot the loss curves with a suitable learning rate for both training data and validation data. What's the issue of period 2 called? Name a possible actions that you could do without changing the number of parameters in your network to counteract this problem.



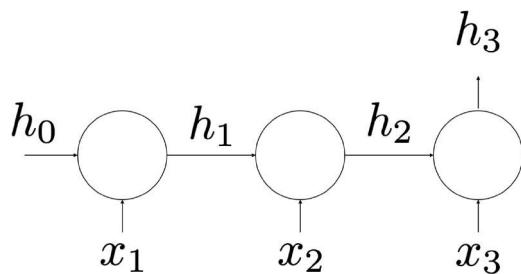
overfitting weight decay
Regulation Mean dev
from anyone drop one





Problem 6 Recurrent Neural Networks and Backpropagation (9 credits)

Consider a vanilla RNN cell of the form $h_t = \tanh(V \cdot h_{t-1} + W \cdot x_t + b)$. The figure below shows the input sequence x_1 , x_2 , and x_3 .



- 0 a) Given the dimensions $x_t \in \mathbb{R}^3$ and $h_t \in \mathbb{R}^5$, what is the number of parameters in the RNN cell? (Calculate final number)

1

- 0 b) If x_t and b are the 0 vector, then $h_t = h_{t-1}$ for any value of h_t . Discuss whether this statement is correct.

1

Now consider the following **one-dimensional** ReLU-RNN cell without bias b .

$$h_t = \text{ReLU}(V \cdot h_{t-1} + W \cdot x_t)$$

(Hidden state, input, and weights are scalars)

- 0 c) Calculate h_2 and h_3 where

$$V = -3, \quad W = 3, \quad h_0 = 0, \quad x_1 = 2, \quad x_2 = 3 \quad \text{and} \quad x_3 = 1.$$

1



$$d \quad 5x3 + 5x5 = 40 = 40$$

b. No, because tanh and V.

$$h_t = \tanh(V \cdot h_{t-1})$$

$$c. \quad h_1 = \text{ReLU}(V \cdot h_0 + W \cdot x_1)$$

$$= \text{ReLU}((-3) \cdot 0 + 3 \cdot 2)$$

$$= \text{ReLU}(6)$$

$$= 6$$

$$h_2 = \text{ReLU}(V \cdot h_1 + W \cdot x_2)$$

$$= \text{ReLU}((-3) \cdot 6 + 3 \cdot 3)$$

$$= 0$$

$$h_3 = \text{ReLU}((-3) \cdot 0 + 3 \cdot 1)$$

$$= 3$$



0
1
2
3

d) Calculate the derivatives $\frac{\partial h_3}{\partial V}$, $\frac{\partial h_3}{\partial W}$, and $\frac{\partial h_3}{\partial x_1}$ for the forward pass of the ReLU-RNN where

$$V = -2, \quad W = 1, \quad h_0 = 2, \quad x_1 = 2, \quad x_2 = \frac{3}{2} \quad \text{and} \quad x_3 = 4.$$

for the forward outputs

$$h_1 = 0, \quad h_2 = \frac{2}{3}, \quad h_3 = 1.$$

Use that $\left. \frac{\partial}{\partial x} \text{ReLU}(x) \right|_{x=0} = 0$.

$$\begin{aligned} \frac{\partial h_3}{\partial x_1} &= \frac{\partial \text{ReLU}(V \cdot h_2 + W \cdot x_3 + b)}{\partial (V \cdot h_2 + W \cdot x_3 + b)} \cdot \frac{\partial (V \cdot h_2 + W \cdot x_3 + b)}{\partial h_2} \cdot \frac{\partial \text{ReLU}(V \cdot h_1 + W \cdot x_2 + b)}{\partial (V \cdot h_1 + W \cdot x_2 + b)} \\ &= \frac{\partial \text{ReLU}(V \cdot h_1 + W \cdot x_2 + b)}{\partial h_1} \cdot \frac{\partial \text{ReLU}(V \cdot h_0 + W \cdot x_1 + b)}{\partial (V \cdot h_0 + W \cdot x_1 + b)} \cdot \frac{\partial (V \cdot h_0 + W \cdot x_1 + b)}{\partial x_1} \\ &= 1 \cdot V \cdot \frac{2}{3} \cdot V \cdot 0 \cdot W = 0 \end{aligned}$$



$$\frac{\partial h_3}{\partial V} = \underbrace{\frac{\partial \text{ReLU}(V \cdot h_2 + w \cdot x_3 + b)}{\partial (V \cdot h_2 + w \cdot x_3 + b)}}_{\frac{\partial (V \cdot h_2 + w \cdot x_3 + b)}{\partial V}} + \underbrace{\frac{\partial \text{ReLU}(V \cdot h_2 + w \cdot x_3 + b)}{\partial (V \cdot h_2 + w \cdot x_3 + b)}}_{\frac{\partial (V \cdot h_2 + w \cdot x_3 + b)}{\partial h_2}} \cdot \underbrace{\frac{\partial \text{ReLU}(V \cdot h_1 + w \cdot x_2 + b)}{\partial (V \cdot h_1 + w \cdot x_2 + b)}}_{\frac{\partial (V \cdot h_1 + w \cdot x_2 + b)}{\partial V}} + \underbrace{\frac{\partial \text{ReLU}(V \cdot h_2 + w \cdot x_3 + b)}{\partial (V \cdot h_2 + w \cdot x_3 + b)}}_{\frac{\partial (V \cdot h_2 + w \cdot x_3 + b)}{\partial h_2}} \cdot \underbrace{\frac{\partial \text{ReLU}(V \cdot h_1 + w \cdot x_2 + b)}{\partial (V \cdot h_1 + w \cdot x_2 + b)}}_{\frac{\partial (V \cdot h_1 + w \cdot x_2 + b)}{\partial V}} + \underbrace{\frac{\partial \text{ReLU}(V \cdot h_0 + w \cdot x_1 + b)}{\partial (V \cdot h_0 + w \cdot x_1 + b)}}_{\frac{\partial (V \cdot h_0 + w \cdot x_1 + b)}{\partial V}}$$

$$\begin{aligned}
 &= l \cdot h_2 + l \cdot V \cdot \frac{2}{3} \cdot h_1 + l \cdot V \cdot \frac{2}{3} \cdot V \cdot 0 \cdot h_0 \\
 &= l \cdot \frac{2}{3} + l \cdot (-2) \cdot \frac{2}{3} \cdot 0 \\
 &= \underline{\underline{\frac{2}{3}}}
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial h_3}{\partial w} &= \frac{\partial \text{ReLU}(v \cdot h_2 + w \cdot x_3 + b)}{\partial (v \cdot h_2 + w \cdot x_3 + b)} \cdot \frac{\partial (v \cdot h_2 + w \cdot x_3 + b)}{\partial w} \\
 &\quad + \frac{\partial \text{ReLU}(v \cdot h_2 + w \cdot x_3 + b)}{\partial (v \cdot h_2 + w \cdot x_3 + b)} \cdot \frac{\partial (v \cdot h_2 + w \cdot x_3 + b)}{\partial h_2} \cdot \frac{\partial \text{ReLU}(v \cdot h_1 + w \cdot x_2 + b)}{\partial (v \cdot h_1 + w \cdot x_2 + b)} \cdot \frac{\partial (v \cdot h_1 + w \cdot x_2 + b)}{\partial x_2} \\
 &\quad + \frac{\partial \text{ReLU}(v \cdot h_2 + w \cdot x_3 + b)}{\partial (v \cdot h_2 + w \cdot x_3 + b)} \cdot \frac{\partial (v \cdot h_2 + w \cdot x_3 + b)}{\partial h_2} \cdot \frac{\partial \text{ReLU}(v \cdot h_1 + w \cdot x_2 + b)}{\partial (v \cdot h_1 + w \cdot x_2 + b)} \cdot \frac{\partial (v \cdot h_1 + w \cdot x_2 + b)}{\partial h_1} \cdot \frac{\partial \text{ReLU}(v \cdot h_0 + w \cdot x_1 + b)}{\partial (v \cdot h_0 + w \cdot x_1 + b)} \\
 &= 1 \cdot x_3 + 1 \cdot v \cdot \underbrace{\frac{2}{3}}_{\partial h_2 / \partial h_2} \cdot x_2 + 1 \cdot v \cdot \underbrace{\frac{2}{3}}_{\partial h_1 / \partial h_1} \cdot v \cdot 0 \cdot w
 \end{aligned}$$

$$= 1 \cdot 4 + 1 \cdot$$



0 e) A Long-Short Term Memory (LSTM) unit is defined as
1 2

$$\begin{array}{ll} \text{o output} & g_1 = \sigma(W_1 \cdot x_t + U_1 \cdot h_{t-1}), \\ \text{forget} & g_2 = \sigma(W_2 \cdot x_t + U_2 \cdot h_{t-1}), \\ \text{update} & g_3 = \sigma(W_3 \cdot x_t + U_3 \cdot h_{t-1}), \\ & \tilde{c}_t = \tanh(W_c \cdot x_t + u_c \cdot h_{t-1}), \\ & c_t = g_2 \circ c_{t-1} + g_3 \circ \tilde{c}_t, \\ & h_t = g_1 \circ c_t, \\ \text{cell state } & / \text{Memory} \end{array}$$

where g_1 , g_2 , and g_3 are the gates of the LSTM cell.

- 1) Assign these gates correctly to the **forget** f , **update** u , and **output** o gates. (1p)
- 2) What does the value c_t represent in a LSTM? (1p)



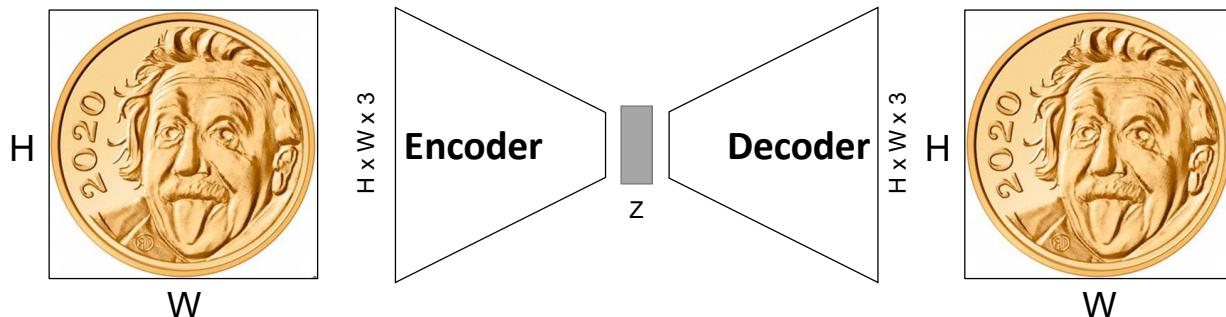


Problem 7 Autoencoder and Network Transfer (11 credits)

You are given a dataset containing 10,000 RGB images with height H and width W of single coins without any labels or additional information.



To work with the image dataset you build an autoencoder as depicted in the figure below:



The input of the encoder is the images of dimension $(H \times W \times 3)$ which are transformed into a one-dimensional real vector with z entries. The latent code is used to decode the input image with the same dimension $(H \times W \times 3)$. Both encoder and decoder are neural networks and the combined network is trainable and uses the L_2 loss as its optimization function.

- a) Is an autoencoder an example of unsupervised learning or supervised learning?

0
1

- b) As the data gets scaled down from the original dimension to a lower-dimensional bottleneck, an autoencoder can be used for data compression. How does an autoencoder as described above differ from linear methods to reduce the dimensionality of the data such as PCA (principal component analysis)?

0
1

- c) For an autoencoder we can vary the size of the bottleneck. Discuss briefly what may happen if

- (i) the latent space is *too small* (1pt).
- (ii) the latent space is *too big* (1pt.)

0
1
2



1. Unsupervised learning

2. Autoencoder is a MVA, PCA is a linear method.
Autoencoder is a MVA with non-linearity
PCA is a linear method

3. If bottleneck too small, Autoencoder may not learn enough features, and will not be able to reconstruct good output

Underfitting

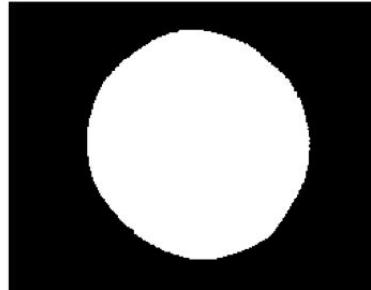
too big, Autoencoder may learn too much they it can remember all the features from training set, include noisy and bias it will poor generalization

Overfitting



- 0 d) Now, you want to generate a random image of a coin. To do so, can you just randomly sample a vector from the latent space to generate a new coin image?

1



- 0 e) Now, someone gives you 1,000 images that are annotated for semantic segmentation of coin and background as shown in the image above. How would you change the architecture of the discussed autoencoder network to perform semantic segmentation?

1

- 0 f) If you wanted to train the new semantic segmentation network what loss function would you use and how?

1

2

- 0 g) How would you leverage your pretrained autoencoder for training a new segmentation network efficiently?

1

2

- 0 h) Why do you expect the pretrained autoencoder variant to generalize more than a randomly initialized network?

1



X

Y

Z

A

B

C

D

E

F

G

H

I

J

K

L

M

N

O

P

Q

R

S

T

U

V

W

X

Y

Z

A

B

C

D

E

F

G

H

I

J

K

L

M

N

O

P

Q

R

S

T

U

V

W

X

Y

Z

A

B

C

D

E

F

G

H

I

J

K

L

M

N

O

P

Q

R

S

T

U

V

W

X

Y

Z

A

B

C

D

E

F

G

H

I

J

K

L

M

N

O

P

Q

R

S

T

U

V

W

X

Y

Z

A

B

C

D

E

F

G

H

I

J

K

L

M

N

O

P

Q

R

S

T

U

V

W

X

Y

Z

A

B

C

D

E

IN-I2DL-1-20200811-E0134-18

d. Autoencoder may not have ability to rebuild a graph from randomly vector.

It's better to use VAE
VAE can sample from the Gaussian latent space to rebuild the com/A

e. convert all the FC layer into fully connected layer
e. Replace the last layer of AE to output 1 or 2 channels

f. Add conv layer with 1 or 2 dim

f. BCE on pixel level / over channels

g.
g. Use pretrain Encoder and frozen weight

h. Access to much more data



Problem 8 Unsorted Short Questions (11 credits)

a) Why do we need activation functions in our neural networks?



b) You are solving the binary classification task of classifying images as cars vs. persons. You design a CNN with a single output neuron. Let the output of this neuron be z . The final output of your network, \hat{y} is given by:

$$\hat{y} = \sigma(\text{ReLU}(z)),$$

where σ denotes the sigmoid function. You classify all inputs with a final value $\hat{y} \geq 0.5$ as car images. What problem are you going to encounter?



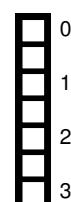
c) Suggest a method to solve exploding gradients when training fully-connected neural networks.



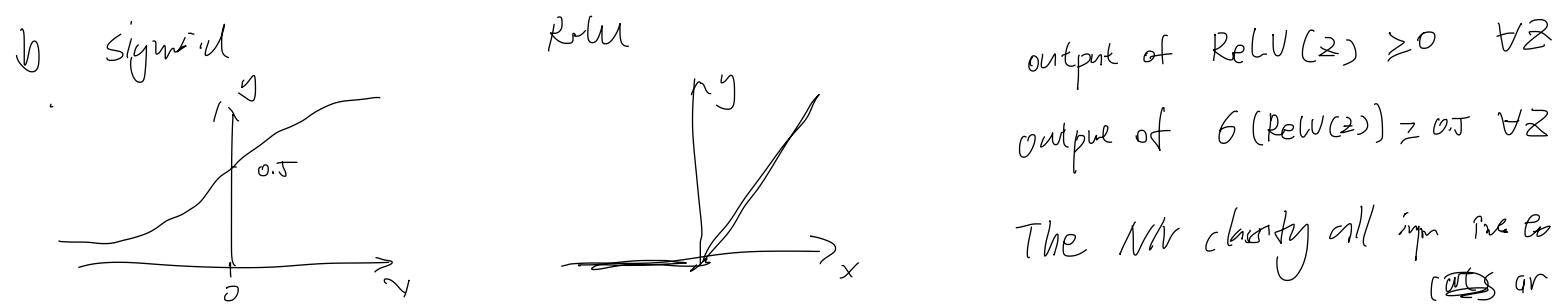
d) Why is it a problem to use the classification accuracy as a loss to train a neural network?



e) Why do we often refer to L_2 -regularization as “weight decay”? Derive a mathematical expression that includes the weights W , the learning rate η , and the L_2 -regularization hyperparameter λ to explain your point.



a. Add non-linearity to the MR so that MR can deal with non-linear data



c. Weight Initialization.

Weight decay Regularization

Gradient clipping

Accuracy

is discrete function
→ no gradient backward ?

d. can not produce the gradient clipping

$$e. w^+ = w - \eta \nabla \left(L(w) + \frac{1}{2} \lambda \|w\|^2 \right)$$

$$w^+ = w - \eta \nabla L(w) - \eta \lambda w$$

$$w^+ = (1 - \eta \lambda) w - \eta \nabla L(w)$$

$$\eta \lambda \ll 1$$

$$\therefore w^+ \rightarrow 0$$



- 0 f) You are given input samples $\mathbf{x} = (x_1, \dots, x_n)$ for which each component x_j is drawn from a distribution with zero mean. For an input vector \mathbf{x} the output $\mathbf{s} = (s_1, \dots, s_n)$ is given by

1

2

3

4

$$s_i = \sum_{j=1}^n w_{ij} \cdot x_j,$$

where your weights w are initialized by a uniform random distribution $U(-\alpha, \alpha)$.

How do you have to choose α such that the variance of the input data and the output is identical, hence $\text{Var}(s) = \text{Var}(x)$?

Hints: For two statistically independent variables X and Y holds:

$$\text{Var}(X + Y) = [\text{E}(X)]^2 \text{Var}(Y) + [\text{E}(Y)]^2 \text{Var}(X) + \text{Var}(X)\text{Var}(Y)$$

Furthermore the PDF of an uniform distribution $U(a, b)$ is

$$f(x) = \begin{cases} \frac{1}{b-a} & \text{for } x \in [a, b] \\ 0 & \text{otherwise.} \end{cases}$$

The variance of a continuous distribution is calculated as

$$\text{Var}(X) = \int_R x^2 f(x) dx - \mu^2,$$

where μ is the expected value of X .



$$\text{Var}(s_i) = \text{Var}\left(\sum_j^N w_{ij} \cdot x_j\right) = \sum_j^N \text{Var}(w_{ij}) \cdot \text{Var}(x_j)$$

$$= N \cdot \left[E(w)^2 \text{Var}(x) + E(x)^2 \text{Var}(w) + \text{Var}(w) \cdot \text{Var}(x) \right]$$

$$X \text{ zero mean} \Rightarrow E(X) = 0$$

$$w \sim U(-\alpha, \alpha)$$

$$\begin{aligned} E(w) &= \int_a^\infty w \cdot f(w) \cdot dw \\ &= \int_a^b \frac{w}{b-a} \cdot dw \\ &= \frac{\frac{1}{2}w^2}{b-a} \Big|_a^b \\ &= \frac{\frac{1}{2}b^2 - \frac{1}{2}a^2}{b-a} \\ &= \frac{b+a}{2} \\ &= \frac{-\alpha + \alpha}{2} = 0 \end{aligned}$$

$$\begin{aligned} \text{Var}(w) &= \int_{-\infty}^{\infty} w^2 \cdot f(w) \cdot dw - E(X)^2 \\ &= \int_a^b \frac{w^2}{b-a} \cdot dw - E(X)^2 \\ &= \frac{1}{3} \cdot \frac{w^3}{b-a} \Big|_a^b - \frac{(a+b)^2}{4} \\ &= \frac{1}{3} \cdot \frac{b^3 - a^3}{(b-a)} - \frac{(a+b)^2}{4} \\ &= \frac{b^2 + ab + a^2}{3(b-a)} \quad (\cancel{b-a}) \\ &= \frac{4b^2 + 4ab + 4a^2 - 3a^2 - 6ab - 3b^2}{12} \end{aligned}$$

$$= \frac{b^2 - ab + a^2}{12}$$

$$= \frac{(b-a)^2}{12}$$

$$= \frac{(\alpha + \alpha)^2}{12} = \frac{\alpha^2}{3}$$

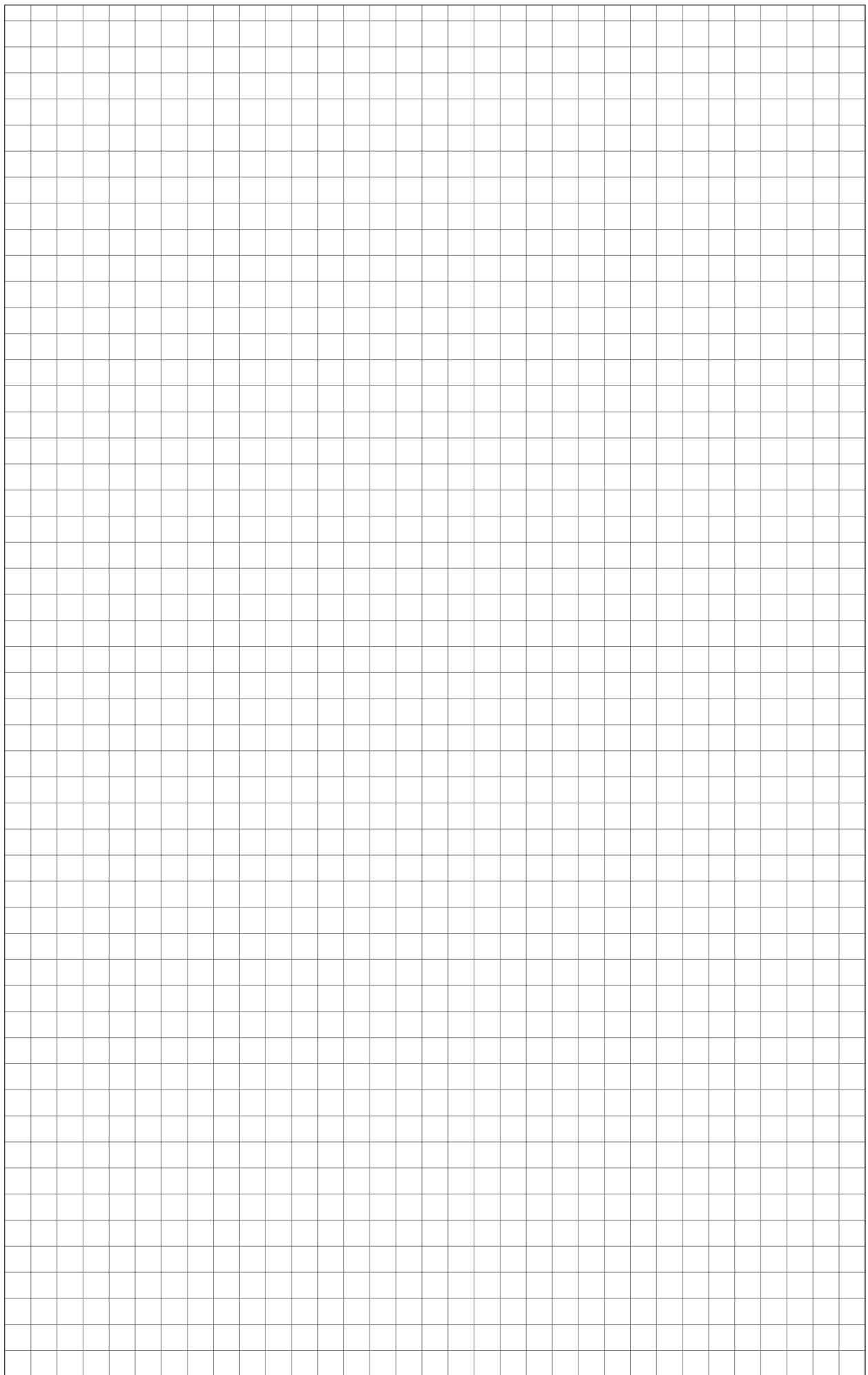
$$n \cdot \frac{\alpha^2}{3} \cdot \text{Var}(x) = \text{Var}(x)$$

$$\alpha = \sqrt{\frac{1}{n}}$$



Additional space for solutions—clearly mark the (sub)problem your answers are related to and strike out invalid solutions.





IN-I2DL-1-20200811-E0134-22

IN-I2DL-1-20200811-E0134-22

IN-I2DL-1-20200811-E0134-22





Y

Y

Y

Y

Y

Y

Y



IN-I2DL-1-20200811-E0134-23





IN-I2DL-1-20200811-E0134-24

IN-I2DL-1-20200811-E0134-24

