

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

Date: Thursday 8th August, 2019

Examiner: Prof. Dr. Leal-Taixé, Prof. Dr. Nießner

Time: 08:00 – 09:30

P 1

P 2

P 3

P 4

P 5

P 6

I						
---	--	--	--	--	--	--

Working instructions

- This exam consists of **20 pages** with a total of **6 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. This includes calculators.

Left room from _____ to _____ / Early submission at _____

6 min

Problem 1 Multiple Choice (18 credits)

Mark your answer clearly by a cross in the corresponding box. Multiple correct answers per question possible.

a) Your network is overfitting. What are good ways to approach this problem?

- Increase the size of the validation set
- Increase the size of the training set
- Reduce your model capacity
- Reduce learning rate and continue training

too much hyperparameters \Leftrightarrow Model can memorize all the training data, then no generalization on new data
high model capacity

b) A sigmoid layer

- has a learnable parameter.
- cannot be used during backpropagation.
- is continuous and differentiable everywhere.
- maps to values between -1 and 1.

↓ ✗

c) Training error does not decrease. What could be a reason?

- Too much regularization.
- Too many weights in your network.
- Bad initialization.
- Learning rate is too high.

too many weights \Rightarrow gradient vanish

best learning ability

d) How many network parameters are in ResNet-152?

- 1,337,337.
- 60,344,232.
- more than a billion.
- 152.

e) 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.

- bcdea
- ebadc
- eadbc
- edbac

e d b a c

f) Dropout



- has trouble with tanh activations.
- is an efficient way for regularization.
- can be seen as an ensemble of networks.
- makes your network train faster.

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

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

h) You are building a model to predict the presence (labeled 1) or absence (labeled 0) of a tumor in a brain scan. The goal is to ultimately deploy the model to help doctors in hospitals. Which of these two metrics would you choose to use?

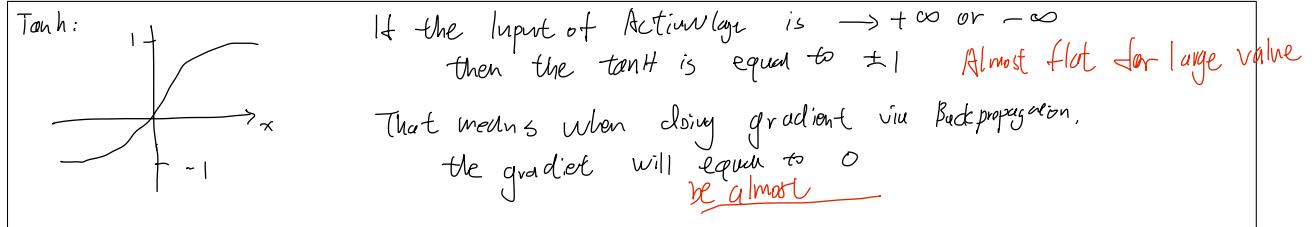
- 真阳性 ← 肿瘤
查漏 ← 没有肿瘤*
- Recall = $\frac{\text{True positive examples}}{\text{Total positive examples}}$. $\frac{\text{TP}}{\text{TP} + \text{FN}}$ 在实际为正样本中预测为正样本的频率
 - Precision = $\frac{\text{True positive examples}}{\text{Total predicted positive examples}}$. $\frac{\text{TP}}{\text{TP} + \text{FP}}$ 在预测为正样本中实际为正样本的频率
 - Average Precision = $\frac{\text{True positive examples} + \text{True negative examples}}{\text{Total examples}}$. 预测模型的误差率

i) Why would you want use 1×1 convolutions? (check all that apply)

- Predict binary class probabilities.
- Collapse number of channels.
- Learn more complex functions by introducing additional non-linearities.
- To enforce a fixed size output.

Problem 2 Short Questions (24 credits)

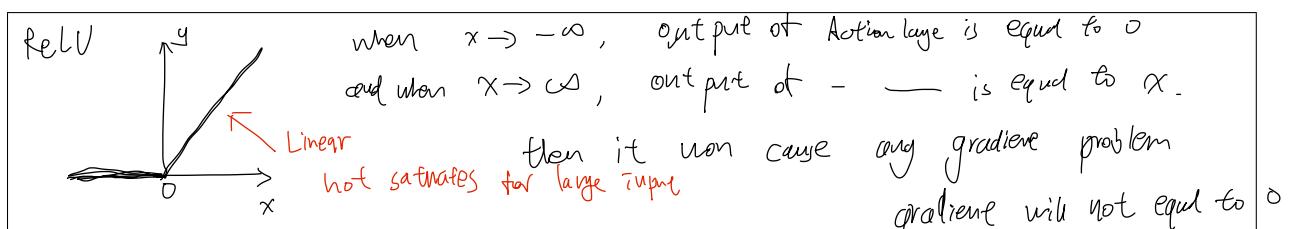
- 0 a) You are training a neural network with 15 fully-connected layers with a tanh nonlinearity. Explain the behavior of the gradient of the non-linearity with respect to very large positive inputs.



- 0 b) Why might this be a problem for training neural networks? Name and explain this phenomenon.

This situation called Gradient ~~saturated~~ / Gradient Vanishing
the weight of MN ~~won't~~ update because of the gradient is ~~equal~~ to 0
very slow

- 0 c) In modern architectures, another type of non-linearity is commonly used. Draw and name this non-linearity (1p) and explain why it helps solve the problem mentioned in the previous two questions (1p).



- 0 d) Why do we often refer to L2-regularization as "weight decay"? Derive a mathematical expression that includes the weights W , the learning rate η , and the L2-regularization hyperparameter λ to explain your point.

Weight update with Loss Function L incl weight decay
 $W = W - \eta \nabla_w (L + \frac{1}{2} \lambda \sum_i \|w_i\|_2^2)$

$$W_{i+1} = (1 - \lambda) W_i - \alpha \nabla L \quad W = W(1 - \underline{\eta \lambda}) - \eta \nabla_w L$$

$$\eta \lambda \ll 1$$

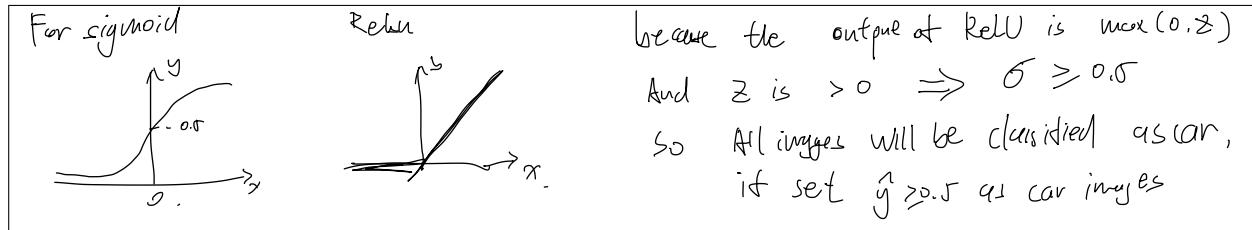
so the W is pushed toward to zero in each iteration

$$\begin{aligned} W^+ &= W - \eta \nabla [L(w) + \frac{1}{2} \lambda \|w\|_2^2] \\ &= W - \eta \nabla L(w) - \eta \lambda w \\ &= (1 - \eta \lambda) w - \eta \nabla L(w) \\ \eta \lambda \ll 1 &\quad \text{Therefore with weight update} \\ &\quad W \text{ will be forced to equal to } 0 \end{aligned}$$

e) 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:

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

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



f) Suppose you initialize your weights w with uniform random distribution $U(-\alpha, \alpha)$. The output s for given input vector x is given by

$$s_i = \sum_{j=0}^n w_{ij} \cdot x_j,$$

where n is the number of input values.

Assume that the input data x and weights are independent and identically distributed. How do you have to choose α such that the variance of the input data and the output is identical, hence $\text{Var}(s) = \text{Var}(x)$.

Hint: 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

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

where μ is the expected value of X .

$$\begin{aligned} \text{Var}(s_i) &= \text{Var}(Wx) = [\text{E}(W)]^2 \text{Var}(x) + [\text{E}(x)]^2 \text{Var}(W) + \text{Var}(W)\text{Var}(x) \\ &= \text{Var}\left(\sum_{j=0}^n w_{ij} \cdot x_j\right) \quad W = \sim U(a, b) = \begin{cases} \frac{1}{2\alpha} & x \in [-\alpha, \alpha] \\ 0 & \text{otherwise.} \end{cases} \quad \text{E}(x) = \int_{-\alpha}^{\alpha} x f(x) dx \\ &= n \text{var}(W \cdot x) \quad = \int_a^b x \cdot \frac{1}{b-a} dx \\ &= \int_{-\alpha}^{\alpha} x \cdot \frac{1}{2\alpha} dx - \mu^2 \quad = \frac{x^2}{2(b-a)} = \frac{b^2 - a^2}{2(b-a)} = \frac{b+a}{2} = \frac{-\alpha+\alpha}{2} = 0 \\ &= \frac{1}{4} \alpha^2 \Big|_0^b \quad \text{Var}(W) = \int_{\mathbb{R}} \frac{\alpha}{2} dx - \mu^2 \\ &= \frac{1}{4} \alpha^2 - \mu^2 \quad = \frac{b^2 - a^2}{2(b-a)} = \frac{b+a}{2} = \frac{-\alpha+\alpha}{2} = 0 \\ &= \frac{1}{2} \frac{(b-a)}{b-a} \quad \text{Var}(x) = \int_{-\alpha}^{\alpha} x^2 f(x) dx - \mu^2 \\ &= \frac{1}{2} (b-a) \quad = \int_a^b \frac{x^2}{3(b-a)} dx - \left(\frac{b+a}{2}\right)^2 \\ &\quad \quad \quad = \frac{b^3 - a^3}{3(b-a)} - \frac{(b+a)^2}{4} \\ &\quad \quad \quad = \frac{(b-a)^2}{12} = \frac{1}{3} \alpha^2 \end{aligned}$$

g) Consider 2 different models for image classification of the MNIST data set. The models are: (i) a 3 layer perceptron, (ii) LeNet.

Which of the two models is more robust to translation of the digits in the images? Give a short explanation why.

$$\text{Var}(x) = \frac{4b^2 + 4a^2 + 4a^2}{12} - \frac{3b^2 + 3a^2 + 3a^2}{12} \quad \text{Var}(w) = \frac{\alpha^2}{3}$$

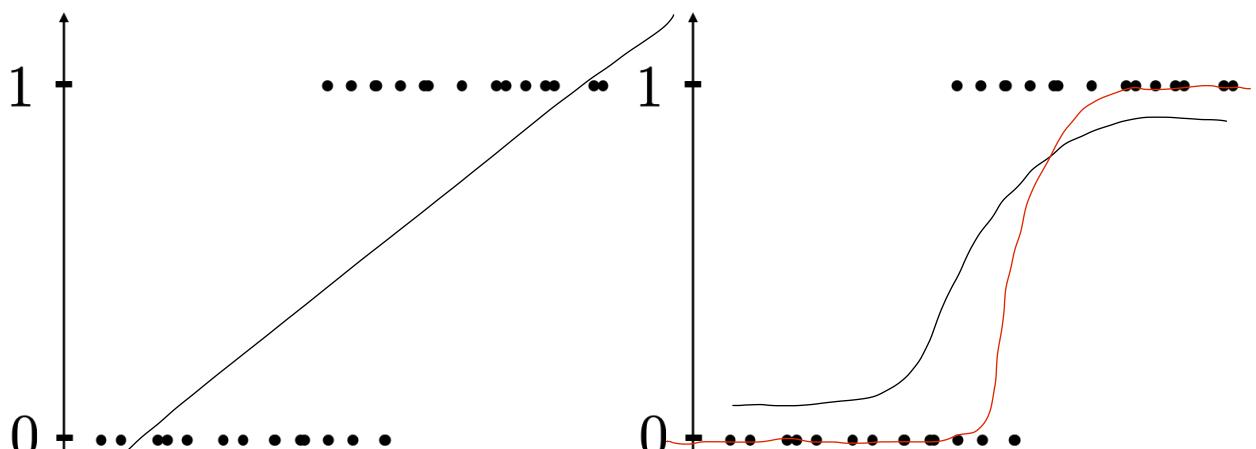
$$= \frac{(b-a)^2}{12}$$

$$n \cdot \frac{\alpha^2}{3} = 1$$

$$\alpha = \sqrt{\frac{3}{n}}$$

LeNet, convolution layers

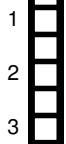
- 0 h) Consider the following one-dimensional data points with classes {0, 1}. Sketch a linear (0.5p) and logistic (0.5p) regression into the figures. Which model is more suitable for this task (1p)?



Plot linear regression (left) and logistic regression (right).

Regression because the label is {0, 1}

- 0 i) What is the mean and standard deviation of Xavier initialization? What changes to this initialization would you propose when used with ReLU non-linearities?



$$\text{Var}(x) = \frac{1}{N} \quad \text{Mean} = 0$$

Use kaiming -Initialization

$$\text{Var}(x) = \frac{2}{N}$$

NN prefer Cat because imbalance between classes

- 0 j) You have 4000 cat and 100 dog images and want to train a neural network on these images to do binary classification. What problems do you foresee with this dataset distribution? Name two possible solutions.



The minibatch may have different distribution, \Rightarrow the NN will cause

Training error not decreasing

collect more dog picture

reweight data loader

Data Augmentation for dog picture / Dropout

- 0 k) Why is initializing all the weights of a fully connected layer to the same value problematic during training?



All neurons have same weights \Rightarrow the gradient via backpropagation on each neuron have same operation \Rightarrow the gradient are same

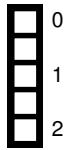
\Rightarrow ~~Learn nothing~~

Node learn the same thing during backpropagation

- I) What is the difference between dropout for convolutional layers compared to dropout for fully connected layers? Explain both behaviours.

Dropout in FCL means use only half ($p=0.5$) neurons to train model

Conv drop feature map random / Fully connected drop weights are random



30 min

If high-variance of gradient, then low learning rate.
low-variance of gradient, then high learning rate.
self-adaptive learning rate method

Problem 3 Optimization (12 credits)

0 a) Explain the concept behind RMSProp optimization. How does it help converging faster?

RMSProp uses second-order Momentum, which called Look-ahead-momentum.
It will predict the location which the optimum will arrive and jump to that location and do some correction.

0 b) Which SGD variation uses first and second momentum?

Adam

0 c) Why is it common to use a learning rate decay? (c) higher gradient
far away closer overshooting

We want fast convergence when it's far away from global optimum so that we can speed up the training process. At mean time, we want slow convergence when it's near to the global opt so that we won't miss or jump over the global optimum point.

0 d) What is a saddle point? What is the advantage/disadvantage of Stochastic Gradient Descent (SGD) in dealing with saddle points?

The goal, the global optimum Saddle, the gradient is zero neither local minima nor local maxima SGD uses stochastic weight update, can help escape saddle point Disadvantage: SGD will take long time to get to saddle point because of stochastic

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

Large mini-batches means in every iteration, SGD will learn more than and SGD will make a better convergence + More efficient computation / better generalization Make gradient less noisy - sensitive to the weight initialization

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

We have too many data. it's very time and hardware costly when doing training on large batches Limited GPU Memory, faster compute, faster update

0 g) Your network's training curve diverges (assuming data loading is correct). Name one way to address the problem through hyperparameter change.

Use small learning rate
reduce

h) What is an epoch?

One

Epoch means the optimization update the weights when using all Minibatch
throughout all training / mini batch in one Epoch
full run through entire training set

i) When is SGD guaranteed to converge to a local minima (provide formula)?

$$\sum_{i=1}^{\infty} \alpha_i = \infty \quad \text{and} \quad \sum_{i=1}^{\infty} \alpha_i^2 < \infty \quad \alpha_i \geq 0 \quad \forall i \geq 0$$

0
1

0
1
2

Problem 4 Convolutional Neural Networks and Advanced Architectures (12 credits)

In the following we assume that the input of our network is a $224 \times 224 \times 3$ color (RGB) image. The task is to perform image classification on 1000 classes. You design a network with the following structure [CONV - RELU] $\times 20$ - FC - FC. That is, you place 20 consecutive convolutional layers (including non-linear activations), followed by two fully-connected layers. Each layer will have its own number of filters and kernel size.

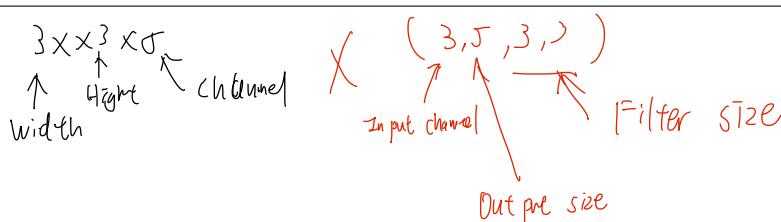
- 0 a) The first 3 convolutional layers have each 5 filters with kernels of size 3×3 , applied with stride 1 and no padding. How large is the receptive field of a feature after the 3 convolutional operations?

$$\begin{aligned} \text{Input} &= F_1 = 5 \\ &\frac{F_1 - 1}{5} + 1 = F_2 & F_2 &= 5 \\ &\frac{F_2 - 1}{5} + 1 = F_3 & F_3 &= 7 \quad 7 \times 7 \\ &\frac{F_3 - 1}{5} + 1 = 1 \end{aligned}$$

- 0 b) What are the dimensions of the feature map after the 3 convolutional operations from (a) ?

$$\begin{aligned} \frac{F_1 - 1}{5} + 1 &= F_2 & \frac{224 - 3}{1} + 1 &= 222 & 218 \times 218 \times 5 \\ && \frac{222 - 3}{1} + 1 &= 220 & \\ && \frac{220 - 3}{1} + 1 &= 218 & \end{aligned}$$

- 0 c) What are the dimensions of the weight tensor of the first convolutional layer? (1p) What does each dimension represent? (1p)



- 0 d) After the 10th convolutional layer your feature map has size $100 \times 100 \times 224$. You realize the next convolutional filter operation will involve too many multiplications that make your network training slow. However, the next layer requires identical spatial size of the feature map.

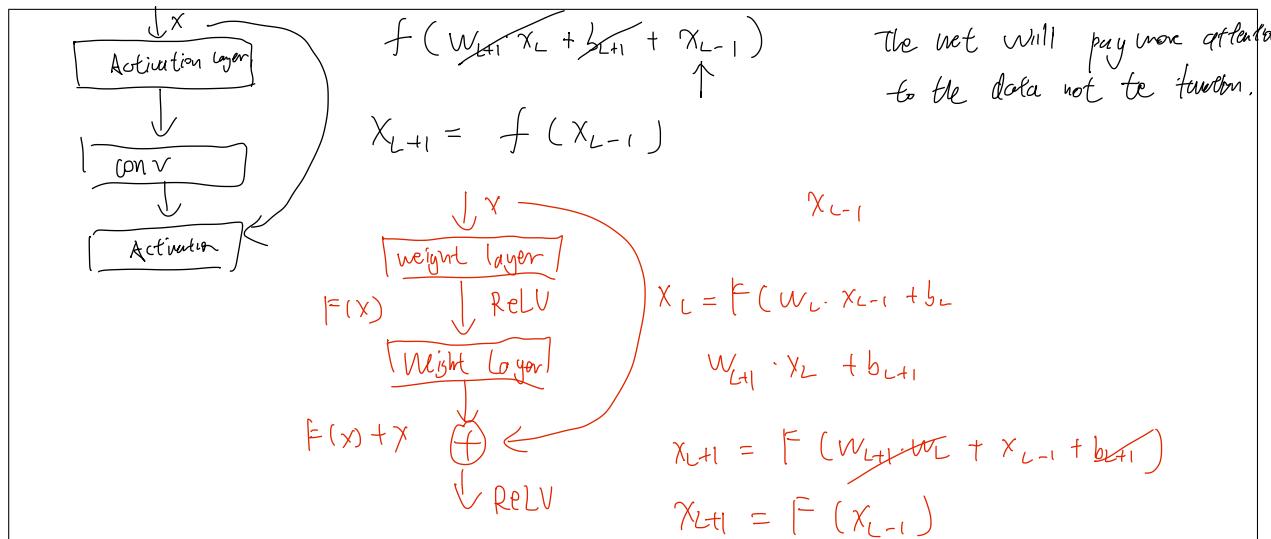
Propose a solution for this problem (1p) and demonstrate your solution with an example (1p).

Use Depth wise separable convolution layer Use 1×1 convolutional layer to reduce the channel	Use $1 \times 1 \times 256$, 256 filter then get $100 \times 100 \times 256$ cheaper and faster for next layer
--	---

- 0 e) Your network is now trained for the task of image classification. You now want to use the trained weights of this network for the task of *image segmentation* for which you need a pixel-wise output. Which layers of your original network described above can you *not* reuse for the image segmentation task? (1p) Describe briefly how you would adapt the network for image segmentation given *any input image size*? (1p)

last FC layer, which is softmax because they take fixed input.
 First adapt the input size to the original network input size by adding layer or padding
 converted FC \rightarrow Fully Convolution \oplus up sampling

f) You decide to increase the number of layers substantially and therefore you switch to a ResNet architecture. Draw a ResNet block (1p). Describe all the operations inside the block (1pt). What is the advantage of using such a block in terms of training (1p)?



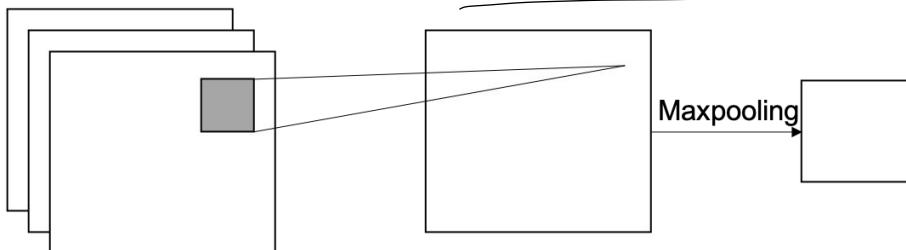
Sum the input and the output from Last Activation Layer

resolving gradient vanishing
highway for gradient / deeper network easier to train

Problem 5 Backpropagation and Convolutional Layers (12 credits)

Your friend is excited to try out those "Convolutional Layers" you were talking about from your lecture. However, he seems to have some issues and requests your help for some theoretical computations on a toy example.

Consider a neural network with a convolutional (without activation) and a max pooling layer. The convolutional layer has a single filter with kernel size $(1, 1)$, no bias, a stride of 1 and no padding. The filter weights are all initialized to a value of 1. The max pooling layer has a kernel size of $(2, 2)$ with stride 2, and 1 zero-padding.



You are given the following input image of dimensions $(3, 2, 2)$:

$$x = \left(\begin{bmatrix} 1 & -0.5 \\ 2 & -2 \end{bmatrix}, \begin{bmatrix} -2 & 1 \\ -1.5 & 1 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \right)$$

0 a) Compute the forward pass of this input and write down your calculations.

Four p

$$\begin{bmatrix} 1 & -0.5 \\ 2 & -2 \end{bmatrix} + \begin{bmatrix} -2 & 1 \\ -1.5 & 1 \end{bmatrix} + \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 0.5 \\ 0.5 & -1 \end{bmatrix}$$

Max pooling

$$\begin{bmatrix} 0 & 0.5 \\ 0.5 & -1 \end{bmatrix} = \begin{bmatrix} 0 & 0.5 \\ 0.5 & 0 \end{bmatrix}$$

1 b) Consider the corresponding ground truth,

$$y = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

Calculate the binary cross-entropy with respect to the natural logarithm by summing over all output pixels of the forward pass computed in (a). You may assume $\log(0) \approx -10^9$. (Write down the **equation** and keep the logarithm for the final result.)

$$\text{BCE : } L = -\frac{1}{n} \sum_{i=1}^n \left[(y_i \cdot \log y_i + (1-y_i) \cdot \log (1-y_i)) \right]$$

$$= - \sum_i t_i \cdot \log y_i \quad \stackrel{\approx}{=} - \sum_i (y_i \cdot \log (w_1 \cdot x_i))$$

$$= - (1 \cdot \log 0.5 + 1 \cdot \log 0.5)$$

$$= -2 \log 2$$

0 c) You don't recall learning the formula for backpropagation through convolutional layers but those 1×1 convolutions seem suspicious. Write down the name of a common layer that is able to produce the same result as the convolutional layer used above.

Fully connected Layer

d) Update the kernel weights accordingly by using gradient descent with a learning rate of 1. (Write down your calculations!)

$$\frac{\partial \text{BCE}}{\partial w_1} = -\frac{\log(w_1 \cdot 2 + w_2 \cdot (-0.5)) + \log(w_1 \cdot (-0.5) + w_2 \cdot 1)}{\partial w_1}$$

$$= -\frac{2}{w_1 \cdot 2 + w_2 \cdot (-0.5)} - \frac{(-0.5)}{w_1 \cdot (-0.5) + w_2 \cdot 1} = -4 + 1 = -3$$

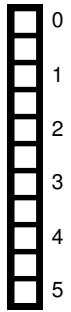
$$\frac{\partial \text{BCE}}{\partial w_2} = -\frac{(-0.5)}{w_1 \cdot 2 + w_2 \cdot (-0.5)} - \frac{1}{w_1 \cdot (-0.5) + w_2 \cdot 1} = 3 - 2 = 1$$

$$w_1 = w_1 - \alpha \cdot \frac{\partial \text{BCE}}{\partial w_1} = 1 - 1 \cdot (-3) = 4$$

$$w_2 = w_2 - \alpha \cdot \frac{\partial \text{BCE}}{\partial w_2} = 1 - 1 \cdot 1 = 0 \quad (3, 1, 1, 1) \Rightarrow$$

$$\frac{\partial \text{BCE}}{\partial w_3} = 0$$

$$w_3 = w_3 - \alpha \cdot \frac{\partial \text{BCE}}{\partial w_3} = 1 - 1 \cdot 0 = 1$$



- 0  e) After helping your friend debugging, you want to showcase the power of convolutional layers. Deduce what kind of 3×3 convolutional filter was used to generate the output (right) of the grayscale image (left) and write down its 3×3 values.

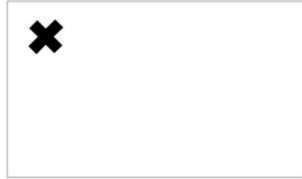
1
2



Edge detection	Vertical edge detector	Prewitt kernel
$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & 8 \\ -1 & -1 & -1 \end{bmatrix}$	$\begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix}$	$\begin{bmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{bmatrix}$
	$\begin{bmatrix} -1 & 0 & 1 \\ -1 & 6 & 1 \\ -1 & 0 & 1 \end{bmatrix}$	$\begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix}$

- 0  f) He finally introduces you to his real problem. He wants to find 3×3 black crosses in grayscale images, i.e., each pixel has a value between 0 (black) and 1 (white).

1



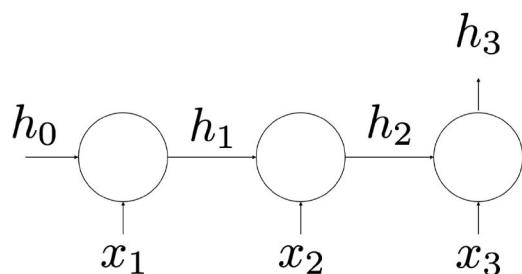
You notice that you can actually hand-craft such a filter. Write down the numerical values of a 3×3 filter that maximally highlights on the position of black crosses.

$$\begin{bmatrix} 2 & 0 & 2 \\ 0 & 1 & 0 \\ 2 & 0 & 2 \end{bmatrix} \quad \begin{bmatrix} -1 & 1 & -1 \\ 1 & -1 & 1 \\ -1 & 1 & -1 \end{bmatrix}$$

Problem 6 Recurrent Neural Networks and LSTMs (12 credits)

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

0
1
2



Given the dimensions $x_t \in \mathbb{R}^4$ and $h_t \in \mathbb{R}^{12}$, what is the number of parameters in the RNN cell? Neglect the bias parameter.

Output $\rightarrow 12$

$$4 \times 12 + 12 \times 12 = 192$$

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

0
1
2

False
there has maybe V and Non-linearly tanh
so $h_t \neq h_{t-1}$

0
1
2
3

c) Now consider the following **one-dimensional** ReLU-RNN cell.

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

(Hidden state, input, and weights are scalars)

Calculate h_1, h_2 and h_3 where $V = 1, W = 2, h_0 = -3, x_1 = 1, x_2 = 2$ and $x_3 = 0$.

For $t=1$

$$V \cdot h_0 + W \cdot x_1 = 1 \cdot (-3) + 2 \cdot 1 = -1$$

$$\text{ReLU}(-1) = 0$$

$$h_1 = 0$$

For $t=2$

$$V \cdot h_1 + W \cdot x_2 = 1 \cdot 0 + 2 \cdot 2 = 4$$

$$\text{ReLU}(4) = 4$$

$$h_2 = 4$$

For $t=3$

$$V \cdot h_2 + W \cdot x_3 = 1 \cdot 4 + 2 \cdot 0 = 4$$

$$\text{ReLU}(4) = 4$$

$$h_3 = 4$$

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 Cell of (c). Use that $\left. \frac{\partial}{\partial x} \text{ReLU}(x) \right|_{x=0} = 1$.

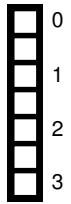
$$h_t = \text{ReLU}(V \cdot h_{t-1} + W \cdot x_t) = \text{ReLU}(z_t)$$

$$\begin{aligned} \frac{\partial h_3}{\partial V} &= \underbrace{\frac{\partial \text{ReLU}(z_3)}{\partial z_3}}_{=1} \cdot \frac{\partial z_3}{\partial V} + \underbrace{\frac{\partial \text{ReLU}(z_2)}{\partial z_2}}_{=0} \cdot \frac{\partial z_2}{\partial V} + \underbrace{\frac{\partial \text{ReLU}(z_1)}{\partial z_1}}_{=0} \cdot \frac{\partial z_1}{\partial V} \\ &= 1 \cdot h_2 + 1 \cdot h_1 + 0 \cdot h_0 \\ &= 4 \end{aligned}$$

$$\begin{aligned} \frac{\partial h_3}{\partial W} &= \underbrace{\frac{\partial \text{ReLU}(z_3)}{\partial z_3}}_{=1} \cdot \frac{\partial z_3}{\partial W} + \underbrace{\frac{\partial \text{ReLU}(z_2)}{\partial z_2}}_{=0} \cdot \frac{\partial z_2}{\partial W} + \underbrace{\frac{\partial \text{ReLU}(z_1)}{\partial z_1}}_{=0} \cdot \frac{\partial z_1}{\partial W} \\ &= 1 \cdot x_3 + 1 \cdot x_2 + 0 \cdot x_1 \\ &= 2 \end{aligned}$$

$$\frac{\partial h_3}{\partial x_1} = \underbrace{\frac{\partial \text{ReLU}(z_3)}{\partial z_3}}_{=1} \cdot \frac{\partial z_3}{\partial h_2} \cdot \underbrace{\frac{\partial h_2}{\partial z_2}}_{=\text{ReLU}(z_2)} \cdot \frac{\partial z_2}{\partial h_1} \cdot \underbrace{\frac{\partial h_1}{\partial z_1}}_{=\text{ReLU}(z_1)} \cdot \frac{\partial z_1}{\partial x_1}$$

$$= 1 \cdot V \cdot 1 \cdot V \cdot 0 \cdot W = 0$$



0
1
2

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

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

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) long-term memory

Cell state

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

A large grid of squares, approximately 20 columns by 25 rows, intended for students to write their solutions. The grid is composed of thin black lines on a white background.

