



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: Tuesday 13th February, 2024

Examiner: Prof. Dr. Matthias Nießner

Time: 10:30 – 12:00

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

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 95 credits.
- Detaching pages from the exam is prohibited.
- **Answers are only accepted if the solution approach is documented.** Give a reason for each answer unless explicitly stated otherwise in the respective subproblem.
- Do not write with red or green colors nor use pencils.
- Do not write outside of the solution boxes.
- Physically turn off all electronic devices, put them into your bag and close the bag.

Left room from _____ to _____ / Early submission at _____

Problem 1 Multiple Choice (18 credits)

Mark correct answers with a cross



To undo a cross, completely fill out the answer option



To re-mark an option, use a human-readable marking



Please note:

- For all multiple choice questions any number of answers, i.e. either zero (!), one 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.**

1.1 Which of the following statements are true about 2D Average Pooling?

- ☒ It does not have any learnable parameters.
- ☐ It always performs better than 2D Max Pooling, since it does not ignore any input features.
- ☐ It reduces the spatial dimensions of the input. Not true. Could be applied with 3, 1, 1
- ☐ It reduces the channel dimension of the input.

1.2 Which of the following statements are true about deep learning optimizers?

- ☐ RMSProp is equivalent to SGD with momentum and learning rate divided by an exponentially weighted gradients average.
- ☐ RMSProp employs a local estimation of the Hessian.
- ☒ RMSProp requires less memory than Adam.
- ☒ Adam optimizer uses 1st and 2nd moment of gradients.

1.3 Which of the following statements are true about transformers?

- ☐ The attention mechanism itself is invariant to order. Not true - it is Equivariant.
- ☐ The concept of transformer can only be applied to text data.
- ☐ Transformers utilize convolutional layers to gather context information.
- ☒ Due to masked attention, the decoder output only depends on the previous outputs and the encoder.

1.4 Which of the following statements are true about autoencoders?

- ☐ The hidden layer should ideally be larger than the input layer.
- ☐ The encoder and the decoder have the same number of layers.
- ☐ Any autoencoder can perfectly learn the identity function.
- ☒ An autoencoder can be used as a lossy compressor.

1.5 Which of the following statements are true about upsampling in neural networks?

- ☐ Applying a convolution followed by a transposed convolution with the identical filter weights is equivalent to an identity operation.
- ☒ Transposed convolution adds trainable parameters to the model, whereas interpolation does not.

Removed

- ☐ Transposed convolution is equivalent to bilinear interpolation (upsampling) followed by applying a standard convolution.

1.6 Given a regression task where the labels are in the range $[-1, 1]$, which of the following final activation functions $f(x)$ could be used?

- ☒ $f(x) = 2\sigma(x) - 1$, where σ corresponds to the sigmoid function.
- ☒ $f(x) = \tanh(10x)$
- ☐ $f(x) = \max(-1, x)$
- ☒ $f(x) = \min(\max(-1, x), 1)$

1.7 Which of the following are true about Generative Adversarial Networks (GANs)?

- ☒ The generator learns by maximizing the probability that a fake image will be classified 'real' by the discriminator.
- ☐ The discriminator aims to learn the distribution of input images, but the generator does not.
- ☐ Removed
- ☐ A fully trained discriminator can be used as a sampling mechanism.

1.8 Given a convolutional layer (kernel size k , padding 1, stride 1, m convolution filters), which of the following might affect the output shape of the convolutional layer (input is a tensor of shape $C \times N \times N$, where C is the number of channels)?

- ☒ N
- ☒ k
- ☒ m
- ☐ C

1.9 Which of the following statements are true about backpropagation?

- ☒ The backpropagated gradient through a tanh activation function is always smaller or equal in magnitude than the upstream gradient.
- ☐ During backpropagation, any of sigmoid/tanh/Leaky ReLU activation functions won't change the sign of gradient. LeakyRelu could have a negative slope alpha
- ☐ Vanishing gradient causes deeper layers to learn more slowly than earlier layers.
- ☐ The derivative of the loss with respect to a specific weight in your network is negative (e.g. -6). That means that decreasing this weight (by a tiny amount) would decrease the loss.

Problem 2 Short Questions (20 credits)

0 ☐
1 ☐
2 ☐

2.1 Consider a binary classification task and two trained classifiers. Classifier A achieves 50% accuracy and classifier B achieves 5% accuracy. With no further training possible, which classifier is more useful (1p) and why (1p)?

(1p) B. (1p) Flipping the output would yield 95% accuracy.
(0.5p) "B has learned something" / if not explaining how B is useful.
(-0.5p) if explaining why A is bad and not why B is useful. 50% accuracy is no better than random guessing, while 5% accuracy implies almost consistent incorrect prediction.

0 ☐
1 ☐
2 ☐
3 ☐

2.2 Explain briefly why positional encoding is used in the transformer architecture, particularly for sequential data such as text (1p). Demonstrate what problem would arise if it were not used (2p).

(1p) Adds information about position of tokens in the sequence.
(1p) (without PE) Attention treats input as unordered set / No inherent processing of order. Only (0.5p) for "model" or "network" instead.
(1p) The order is important for understanding the meaning. / example showing this. **Common mistakes:** "What to pay attention to", vague statements like "more context"

0 ☐
1 ☐
2 ☐
3 ☐

2.3 Are the following two functions appropriate to be activation functions in a backpropagation neural network? If not, why?

• a) $y = \cos(x)$, where $y' = -\sin(x)$

• b) $y = 0.5 \cdot \text{sgn}(x)$, where $\text{sgn}(x) = \begin{cases} 1, & \text{if } x > 0 \\ 0, & \text{if } x = 0 \\ -1, & \text{if } x < 0 \end{cases}$

a) (1p) Yes. For being non-linear. Or No, since it fits very specific scenarios.
b) (1p) No. Gradients are zero almost everywhere. (0.5p) for constant gradient everywhere but not explicitly zero. (0p) for "discontinuity" or "non-differentiable".

2.4 Describe the concept of gradient clipping in training neural networks (2p) and explain why it benefits the training of recurrent neural networks particularly (1p).



(2p) Gradient clipping involves capping the values of gradients during backpropagation to some range, to prevent them from becoming excessively large (or small). (-0.5) for specifying a specific range, not as an example. (1p) exploding gradients
(1p) RNNs are prone to vanishing or exploding gradients due to the sequential data. (0.5) if said that clipping/scaling makes the training more stable.

2.5 During training of a neural network, you noticed that the validation loss is smaller than the training loss. Name two possible explanations for this.



Accepted (each 1p):
Easier / simpler validation set (0p for "validation set is easier to learn".) || Training loss is evaluated during the entire epoch, while validation loss is only at the end of the epoch. || Dropout || Bug / error in / wrong implementation.
Partially correct (0.5p):

Regularization applied differently during training and validation".

2.6 Show mathematically that sequentially applying two linear layers with biases (W_1, b_1) and (W_2, b_2) is equivalent to using one linear layer if no non-linearities are applied between them (1p). What will be the weight and bias of this new layer (2p)? You may use x to denote the input.



$$f_2(f_1(x)) = W_2(W_1x + b_1) + b_2 = W_2W_1x + W_2b_1 + b_2$$

-0.5p for small errors in the derivation.

2.7 Explain the β_1 and β_2 hyperparameters in the Adam optimizer.



β_1 : Decay rate coefficient of the moving average of the first moment (0.5p)
for most recent/past/accumulated (0.5p) (or when mentioned mean (0.5p)) gradients.

β_2 : Decay rate coefficient of the moving average of the second moment (0.5p)
for most recent/past/accumulated squared (0.5p) (or when mentioned variance (0.5p)) gradients.

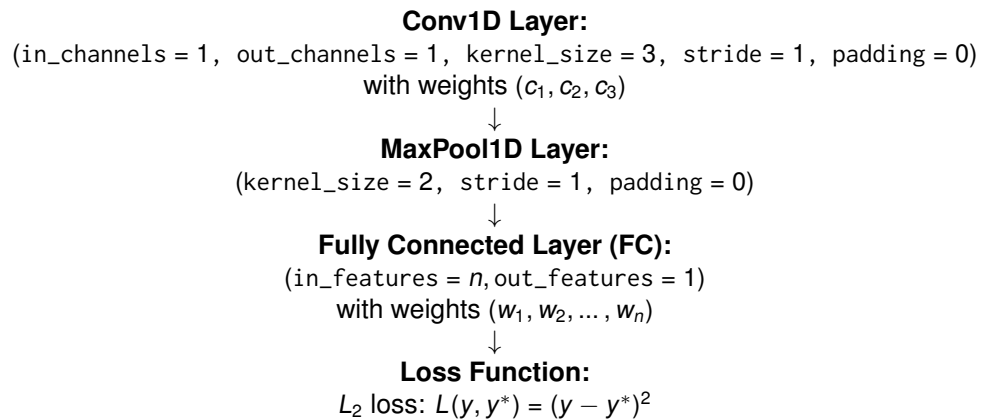
2.8 Explain the Inception layer and list its internal elements.



- idea: multiple sizes conv kernels (0.5p)
- components: convs of different kernels (1, 3, 5), max pool, concat features (1.5p). Also drawings are accepted.

Problem 3 Backpropagation (11 credits)

Consider the following 1D Convolutional Neural Network:



- 0 ☐ 3.1 Consider a **7-dimensional vector**: $\vec{x} = [x_1, x_2, x_3, x_4, x_5, x_6, x_7]$ as input to the network. Write down the output dimensions of the Conv1D and MaxPool1D layers, as well as the input dimension n of the FC layer.

1 ☐

Convolutional Layer Output Dimension: $\frac{n+2p-k}{s} + 1 = \frac{7+0-3}{1} + 1 = 5$
 MaxPool Layer Output Dimension: $\frac{n-k}{s} + 1 = \frac{5-2}{1} + 1 = 4$
 FC Input Dimension: 4
 (1p) for all three answers. (0.5p) for at least 1 correct answer.

For the next questions, we change the network to handle a **5-dimensional input** $\vec{x} = [x_1, x_2, x_3, x_4, x_5]$. Assume the following setting: $\vec{x} = [2, 3, 1, 2, -1]$, convolutional weights $[0.5, 1.5, 0.5]$, FC weights $\vec{w} = [-1, +3]$, and ground-truth label $y^* = 2$.

- 0 ☐ 3.2 Calculate the forward pass through each of the network's layers, including the output y and the loss value $L(y, y^*)$. Suggested notation: \vec{z} for the output of the convolutional network, \vec{m} for the maxpool layer.

1 ☐

2 ☐

3 ☐

$z_1 = 1 + 4.5 + 0.5 = 6$
 $z_2 = 1.5 + 1.5 + 1 = 4$
 $z_3 = 0.5 + 3 - 0.5 = 3$
 $m_1 = \max(6, 4) = 6$
 $m_2 = \max(4, 3) = 4$
 $y = -1 \cdot 6 + 3 \cdot 4 = 6$
 $L = (6 - 2)^2 = 16$
 (1p) for conv output.
 (0.5p) for the maxpool output.
 (1) for the FC input.
 (0.5) for the loss

3.3 Construct a matrix that can be used to express the convolutional layer as a matrix-vector multiplication between the matrix and the input \vec{x} .



$$\begin{bmatrix} z_1 \\ z_2 \\ z_3 \end{bmatrix} = \begin{bmatrix} c_1 & c_2 & c_3 & 0 & 0 \\ 0 & c_1 & c_2 & c_3 & 0 \\ 0 & 0 & c_1 & c_2 & c_3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{bmatrix}$$

(1p) for correct shape: 3×5 or 5×3 .

(1p) for correct values in the matrix, also in a transposed matrix

(-0.5p) for not using zeros in the matrix.

3.4 Given the loss function $L = (y - y^*)^2$, calculate the derivative $\frac{dL}{dy}$. Also, evaluate its explicit value.



$$\frac{dL}{dy} = 2(y - y^*) = 2(6 - 2) = 8$$

(0.5p) derivative || (0.5) evaluation.

Common Mistakes: swapped values for y and y^* , used loss value for y .

0 ☐
1 ☐

3.5 Denote $\vec{m} = [m_1, m_2]$ the outputs of the MaxPool1D layer. Given $\frac{dL}{dy} = 2$, calculate $\frac{dL}{d\vec{m}}$.

$$\frac{dL}{dm_1} = \frac{dL}{dy} \cdot \frac{dy}{dm_1} = 2 \cdot w_1 = 2 \cdot (-1) = -2$$

$$\frac{dL}{dm_2} = \frac{dL}{dy} \cdot \frac{dy}{dm_2} = 2 \cdot w_2 = 2 \cdot 3 = 6$$

(0.5p) for chain rule, (0.5p) for correct solution.

0 ☐
1 ☐
2 ☐

3.6 Denote $\vec{z} = [z_1, z_2, z_3]$ the outputs of the Conv1D layer. Given $\frac{dL}{dm} = [2, 1]$, calculate $\frac{dL}{d\vec{z}}$.

$$\frac{dL}{dz_1} = \frac{dL}{dm} \cdot \frac{dm}{dz_1} = \frac{dL}{dm_1} \cdot \frac{dm_1}{dz_1} + \frac{dL}{dm_2} \cdot \frac{dm_2}{dz_1} = 2 \cdot 1(z_1 > z_2) + 1 \cdot 0 = 2$$

$$\frac{dL}{dz_2} = \frac{dL}{dm} \cdot \frac{dm}{dz_2} = \frac{dL}{dm_1} \cdot \frac{dm_1}{dz_2} + \frac{dL}{dm_2} \cdot \frac{dm_2}{dz_2} = 2 \cdot 1(z_1 < z_2) + 1 \cdot 1(z_2 > z_3) = 1$$

$$\frac{dL}{dz_3} = \frac{dL}{dm} \cdot \frac{dm}{dz_3} = \frac{dL}{dm_1} \cdot \frac{dm_1}{dz_3} + \frac{dL}{dm_2} \cdot \frac{dm_2}{dz_3} = 2 \cdot 0 + 1 \cdot 1(z_2 < z_3) = 0$$

(0.5p) if the chain rule is written but dm/dz is not calculated. || (1.0p) if dm/dz is calculated || (0.5p) if dL/dz is calculated from dL/dm and dm/dz . || -0.5p if solution is given but the approach is not documented.

0 ☐
1 ☐

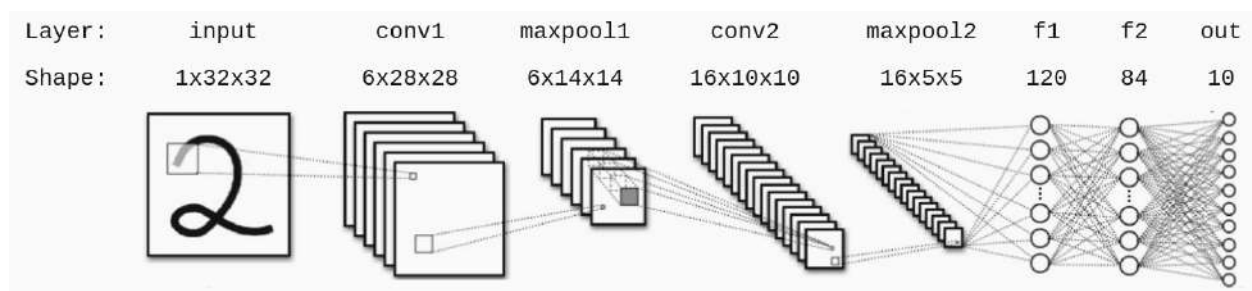
3.7 Given $\frac{dL}{dz} = [1, 2, 3]$, calculate $\frac{dL}{dc_2}$. c_2 is the second convolutional weight.

$$\frac{dL}{dc_2} = \frac{dL}{dz_1} \cdot x_2 + \frac{dL}{dz_2} \cdot x_3 + \frac{dL}{dz_3} \cdot x_4 = \text{value}$$

$$= 1 \cdot 3 + 2 \cdot 1 + 3 \cdot 2 = 11$$

Problem 4 CNN (15 credits)

You are contemplating design choices for a convolutional neural network for the classification of digits. LeCun et. al suggest the following network architecture:



For clarification: the shape **after** having applied the operation 'conv1' (the first convolutional layer in the network) is $6 \times 28 \times 28$. No padding is being used. Convolutions use stride of 1,

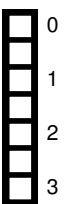
Maxpool uses a stride of 2 and a kernel size of 2.

4.1 Explain the term 'receptive field' (1p). What are the receptive fields of a pixel in the feature maps after the operations 'conv1' and 'maxpool1'? (0.5p each) What is the receptive field of a neuron in layer 'f1' (1p)?

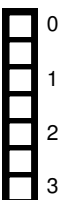
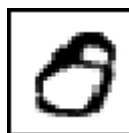
(1p): The receptive field is the size of the region in the input space that a pixel in the output space is affected by. (0.5p) for an incomplete definition.

(0.5p each): Conv: 5x5 (0.5p), Maxpool: 6x6 (0.5p)

(1p): f1: 32x32 (whole image)



4.2 You and your classmates have each trained their own softmax classifier with a multiclass cross-entropy loss. You test your models with the following image of the digit 8, and obtain the outputs as depicted in the table below:



Model \ Digit	0	1	2	3	4	5	6	7	8	9
Model A	0.45	0.00	0.00	0.00	0.05	0.00	0.05	0.05	0.40	0.00
Model B	0.47	0.00	0.00	0.00	0.00	0.00	0.02	0.00	0.41	0.10
Model C	0.30	0.00	0.05	0.00	0.00	0.00	0.20	0.00	0.35	0.10

Table 4.1: Vector of class probabilities, as outputted by models A, B, and C

Which model/models would predict the class correctly for this sample (1p)?

On this sample, which model would yield the largest loss value (1p)? What is its value (1p)? *Note:* You may use elementary functions like log and exp in your answer.

- (1p) Model C.
- (1p) Model C. (0.5p) for not explicitly stating the model, but stating the correct value
- (1p) $-1 * \log(0.35)$, (0.5p) for correct value, w/o the sign / Wrong coefficient (0.1) / wrong value inside a correct loss function (e.g. $-1 * \log(0.41)$)

0 ☐ 4.3 Instead of recognizing MNIST digits, you now want to be able to classify whether an image depicts an even or odd digit. What is a single change to the network architecture needed to support this? Choose a solution that is efficient in terms of number of trainable weights.

(0.5p) Change the number of output neurons from 10 to 1 or 2, or any reasonable method for binary classification.

(0.5p) for efficient solution: i.e. change the number of output neurons from 10 to 1.

0 ☐ 4.4 Instead of taking 32×32 images, you now want to train the network to classify images of size 68×68 . Describe two possible architecture changes to support this.

The most commonly accepted solutions are ((1p) for each):

- Resize layer to downsample images to 32×32 .
- Add an initial convolutional layer with a 5×5 kernel and a subsequent 2×2 max pooling layer before the current architecture (must state the kernel size and stride of the pooling layer).
- Change input dimension of layer f1 to $16 \times 14 \times 14$ (= 3136).
- changing the fully connected layer to 1×1 conv layer only with proper design (e.g. make the final output size $10 \times 1 \times 1$ to have 10 outputs) to accomplish the classification task (0.5P for only mentioning 1×1 conv/fully convolutional network)

(0p) for data augmentation or preprocessing the image (not architecture change) / for autoencoder (latent space commonly 1D space, cannot use encoder to downsample).

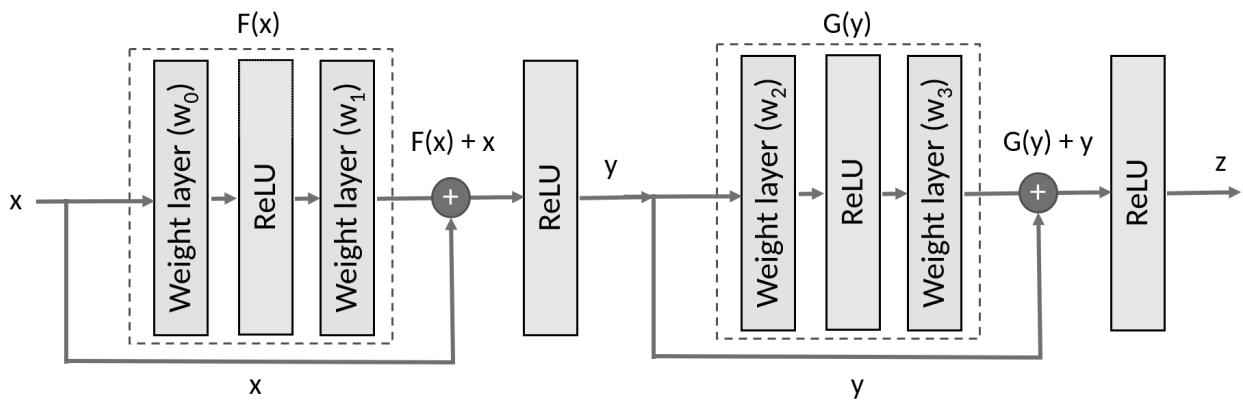
0 ☐ 4.5 Your architecture works and you manage to classify digits fairly well. After reading many online blogs, you decide to try out a much deeper network to boost the network's capacity. Describe two problems that you might encounter when training very deep networks.

(1p) For effect (0.5p) and cause (0.5p) for each problem

- (effect): Vanishing/Exploding Gradients, (cause): too many layers, chain rule, saturation.
- (effect): High Compute/Out of Memory, (cause): too many parameters.
- (effect): Require many steps, (cause): vanishing gradients.
- (effect): Loss diverges/unstable, (cause): exploding gradients
- (effect): Overfitting/Less Generalization, (cause): high network capacity, can memorize
- (effect): Dead ReLUs, (cause): exploding gradients

<input type="checkbox"/>	0
<input type="checkbox"/>	1
<input type="checkbox"/>	2
<input type="checkbox"/>	3
<input type="checkbox"/>	4

4.6 You read that skip connections are beneficial for training deep networks. The following image shows a segment of a very deep architecture that uses skip connections. How are skip connections helpful (1p)? Demonstrate this mathematically by finding the derivative of w_0 (3p). Hint: consider the chain rule, up-stream gradients.



(1p) Explanation: gradient flow highway, which helps to avoid vanishing gradients.

(3p) Math: Assume that $dReLU(x)$ is 1 for all, since it's not interesting when it's 0. What we need to show here, is the highway of gradients. We will do that by showing it in $\frac{dz}{dy}$:

$$\frac{dz}{dw_0} = \frac{dz}{dy} \cdot \frac{dy}{dw_0}$$

$$\frac{dz}{dy} = \frac{(G(y) + y)}{dy} = \left(\frac{dG(y)}{dy} + 1 \right) = w_3 w_2 + 1$$

$$\frac{dy}{dw_0} = \frac{(w_1(w_0 x))}{dw_0} = w_1 x$$

$$\frac{dz}{dw_0} = (w_3 w_2 + 1) w_1 x$$

It is enough to show $\frac{dz}{dy}$, that shows that we have a derivative 1, that represents the highway of gradients (simply keeping the full magnitude of the upstream gradient).

- 1p for correct conclusion, i.e. point out that skip connection introduces the additive term, 0.5p if only verbally explained rather than mathematically.
- 2p if derivative is fully expanded, halved if not expanded (-0.5p for each small mistake).

Problem 5 Convolutional Kernel (7 credits)

0 ☐
1 ☐
2 ☐

5.1 Compute the forward pass for the given input and convolution kernel as follows. Note that in this case we have a stride of 1 and no padding.

$$\text{Input: } \begin{bmatrix} 1 & 2 & -1 & -3 \\ -1 & 3 & 0 & 2 \\ 1 & 1 & 1 & -2 \\ 0 & 2 & 0 & -3 \end{bmatrix}$$
$$\text{Kernel: } \begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$$

Output:

$$\begin{bmatrix} 13 & -5 \\ -2 & 6 \end{bmatrix}$$

(-0.5p) for not simplifying the calculation.

0 ☐
1 ☐

5.2 Actually, the kernel we used above has a particular function. Explain what kind of effect this kernel has when applying it on an image.

Sharpen the image.

Common mistakes: edge detection, corner detection, blurring the image, smoothing the image, increasing the contrast and point/circle detection.

0 ☐
1 ☐
2 ☐
3 ☐
4 ☐

5.3 Compute the gradient with respect to the input, using the above kernel and the following up-stream gradient.

Up-stream gradient:

$$\begin{bmatrix} 3 & -1 \\ 1 & 2 \end{bmatrix}$$

Output:

$$\begin{bmatrix} 0 & -3 & 1 & 0 \\ -3 & 15 & -10 & 1 \\ -1 & 0 & 10 & -2 \\ 0 & -1 & -2 & 0 \end{bmatrix}$$

Full solution will be posted on Piazza. -0.5p for calculation errors. However, points were given for even trying. 0p for assuming that the gradient's shape is 2x2.

If by mistake calculated the gradient with respect to the kernel, also accepted:

$$\begin{bmatrix} 6 & 10 & 4 \\ -3 & 12 & -5 \\ 6 & 4 & -1 \end{bmatrix}$$

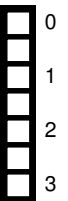
Problem 6 Activation Function (9 credits)

In November 2015, researchers proposed the activation function “Exponential Linear Unit” (ELU):

$$\text{ELU}(x) = \begin{cases} x & \text{if } x \geq 0 \\ \alpha \cdot (e^x - 1) & \text{if } x < 0 \end{cases} \quad (6.1)$$

6.1 Sketch the function $\text{ELU}(x)$ as a graph. Consider the effect of α in your graph. Describe the behavior of $\text{ELU}(x)$ for the following cases:

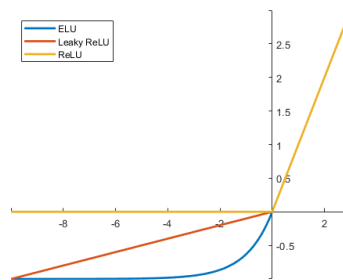
- positive values of x
- large negative values of x
- small negative values of x



The ELU function will behave as follows:

- (1p) for the sketch on the euclidean plane (not a computational graph). (0.5p) for showing the effect of α .
- (0.5p) positive values: identity function (gradient 1, no change in slope at $x = 0$).
- (0.5p) large negative values: close to $-\alpha$
- (0.5p) low negative values: from $-\alpha$ to 0

Common mistakes: not showing the effect of α , Describing the derivative of $\text{ELU}'(x)$.



6.2 Write down the derivative of the ELU function (2p). Express the derivative in terms of the ELU function itself (1p).

$$\text{ELU}'(x) = \begin{cases} 1 & \text{if } x > 0 \\ \alpha e^x & \text{if } x \leq 0 \end{cases}$$



6.3 Describe 2 advantages and 1 disadvantage of $\text{ELU}(x)$ compared to $\text{ReLU}(x)$?

True only for ELU with alpha = 1

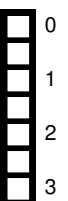
smoothness or continuity: Unlike ReLU, ELU is smooth everywhere, including at $x = 0$, which can help alleviate some optimization issues during training, particularly for deep neural networks.

handles negative inputs: ELU can handle negative inputs gracefully, preventing the "dying ReLU" problem where neurons output zero for negative inputs, which can occur in ReLU when a large fraction of neurons become inactive during training. ELU can produce negative outputs for negative inputs, allowing neurons to be active throughout training.

Disadvantage: more expensive to compute (exponent vs. 0 in the negative part of x)

(0.5p) for points with vague or no description.

(0p) for points that are valid for both ELU and ReLU.



Problem 7 Sparsity (7 credits)

0 ☐ 7.1 Name a regularization method that encourages weights to be sparse. Write down its explicit formula
1 ☐ (assume it is applied on a weights matrix $W_{m \times n}$).

L1 regularization (or Lasso Regression)

$$\|W\|_1 = \sum_{i=1}^m \sum_{j=1}^n |w_{ij}|$$

0 ☐ 7.2 Name one non-learning based approach and one learning based approach to reduce the dimensionality
1 ☐ of data.

Non-learning (0.5p): PCA / SVD / Pooling / Downsampling Layer / Interpolation or Resize.

Learning (0.5p): Autoencoder / Convolution / FC ($n_{output} < n_{input}$).

Incorrect: VAE / Manual feature selection / Cropping / Dropping dimensions.

0 ☐ 7.3 You trained a classification network on 100 classes, which achieves 95% accuracy on your train and 90%
1 ☐ on your validation set. You have deployed your model and noticed that it performs poorly on certain classes.
2 ☐ When examining the class your model performs worst on, you realize that your model only achieves less
than 32% accuracy on that class (both in train and validation). Assuming the data labels are correct, explain
a potential reason, why this can happen and how to solve this problem.

(1p) Class Imbalance / Sparse Dataset || (0.5p) Less / Not enough data (without mentioning minority class)

(1p) Worst performing class comes from a different domain || (0p) Train and Val datasets come from different distributions

Valid Solution: (1p) Augment Minority Class / Stratified Sampling / Loss Weighting / Oversample/Undersample to balance classes || (0.5p) Balance the dataset / Add more data (not answering how) || (0p) Reshuffle

0 ☐ 7.4 You trained a very large model consisting of many linear layers. You have examined the trained model's
1 ☐ weights and noticed many of the weights are very small compared to others in the same layer. Suggest
2 ☐ how sparsity in your weight matrices could be leveraged for a lossy model compression (2p). What are the
3 ☐ practical advantages of doing that (1p)?

(2p) Set Weights to zero below threshold, use sparse matrices

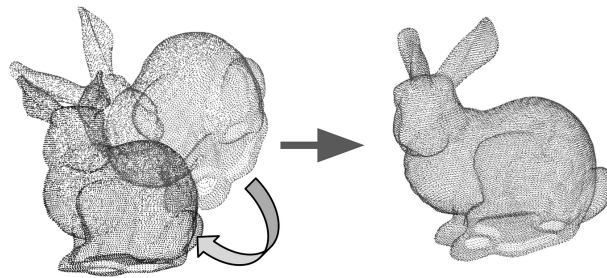
(1p) One of the following:

- More efficient compute

- Less Storage

(0p) Data Compression Feature Extraction

Problem 8 3D Pointclouds (8 credits)



Pointclouds are a simple data structure to represent 3D objects as an unordered set of points.

$$X = \{[x_{i1}, x_{i2}, x_{i3}] \in \mathbb{R}^3 | i = 1, \dots, N\}$$

A batch of pointclouds is then stored in a tensor of shape (B, N, D) , where B - batch size, N - number of points (assume it is the same for all pointclouds), D - dimension (usually 3).

Rigid transformations, like rotation and translation, are common operations on pointclouds. In 3D, a rotation can be represented by a 3×3 rotation matrix (R) and translation by a 1×3 vector (t), similar to an affine layer. The rotated pointcloud X' can be obtained by $X' = X \cdot R + t$.

$$X' = \begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ \vdots & \vdots & \vdots \\ x_{N1} & x_{N2} & x_{N3} \end{bmatrix} \times \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{bmatrix} + \begin{bmatrix} t_1 & t_2 & t_3 \end{bmatrix}$$

Consider the bunny shaped pointcloud as depicted in the image above, and a dataset consisting of copies of the pointcloud, each pointcloud randomly rotated and translated in space.

Your task is to design a neural network to align the pointclouds: given any such pointcloud of the bunny, your network will predict the rotation and translation, that applying them would transform the pointcloud back to its canonical position $(0, 0, 0)$ and orientation.

Assume a supervised setting where the dataset contains ground truth rotation matrices and translation vectors for each pointcloud.

8.1 Network architecture. Explain why a Fully Connected layer (FC) would be a better choice to use on the pointclouds rather than a convolutional layer.

Answer: Point clouds are unordered (1p) (their structure doesn't change with permutations). Therefore, convolutions, that capture features in local neighborhoods, are irrelevant (1p).

Also accepted: (1p): Convs catch Local vs Fc capture global features.

Common mistakes: Convolutions are rotations or translation invariant / variant - irrelevant here. A point cloud is not the same structure as a 2D image. / Convolution would compress the data / FC takes into consideration all points / Has a receptive field of the entire 3D space

0
1
2

8.2 Reshaping input data. How can you reshape the input batch of shape $(B, N, 3)$ so that the first FC layer extracts global information directly from all points of each pointcloud?

Answer: $(B, N \times 3)$ / flatten the second and third dimensions.

Formulation could be confusing, students falsely understood it meant to process the entire batch all together. 0.5p for that (saying the FC input is $B \times N \times 3$ / $(B \times N, 3)$), although it breaks the abstraction of a batch. (0.5p) for "flatten". $(3B, N)$ makes no sense and also breaks the batch abstraction.

0
1

0 ☐ 8.3 Input and output layer sizes. Given a batch size of $B = 8$ pointclouds, with $N = 100$ points each, what are the sizes of the *input layer* and the *output layer* of your fully connected network?

1 ☐

2 ☐ (1p) input: $N * 3 = 100 * 3 = 300$. (1p) output: $9 + 3 = 12$ for both rotation and translation matrices. (4 quaternions and 3 Euler angels are also accepted.). Chain errors from 8.2 were not penalized here. Also accepted: If 8.2 is wrong for misunderstanding, the answer (e.g. 2400 for $B \times N \times 3$) is accepted.

0 ☐ 8.4 Loss function for rotation matrices. Rotation matrices belong to a special group of matrices called Orthogonal Matrices, which observe the following property: its transpose is also its inverse.

1 ☐ Once you have estimated a rotation matrix with your network, you notice that applying R_{pred} to the pointcloud causes skewing to its shape. Your colleague told you it was happening because the network did not output an orthogonal matrix, making R_{pred} an invalid rotation matrix.

2 ☐ Suggest a regularization term $Reg(R_{pred}, R_{GT})$ that would also take into consideration this crucial characteristic of valid rotation matrices and encourage your output matrix to be orthogonal.

3 ☐ Hint: Avoid matrix inversion, because they are costly and not numerically stable.

L1, L2, Frobenius, ... applied like this:

$LossFunc(R_{pred}^T \cdot R_{true}, I)$ or $Loss(R_{pred}^T \cdot R_{true} - I)$.

For anyone who stated any form of matrix operation: -0.5p for not mentioning which reg/loss function they should use, -0.5p for no Identity matrix or 1, -0.5p for false usage of Transpose, -0.5p for not distinguishing between the pred and gt, -0.5p for not performing matrix multiplication between R_{pred} and R_{gt} .

Trying to use the determinant ($==1$) Got 1p for the creative idea, but it is false mathematically (Not every matrix with a $det=1$ is a rotation matrix)

This image shows a full page of blank graph paper. The grid consists of thin, light gray horizontal and vertical lines that intersect to form small squares across the entire surface. There are no margins, text, or other markings on the paper.

