

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 8th February, 2022

Examiner: Prof. Dr. Matthias Nießner

Time: 15:00 – 11:30

Working instructions

Important: The Date and Time stated above refer to the time when the blackened exam is being released. The actual working time of this exam is February 10th, 10:00 - 11:30.

- The blackened exam has the same layout as the non-blackened exam with the acutal questions, which is going to be released once the working time starts.
- Only submit your personalized blackened exam. **DO NOT submit the non-blackened/non-personalized exam** (clearly indicated with "DO NOT SCAN/UPLOAD").
- This final exam consists of **16 pages** with a total of **7 problems**. Please make sure now that you received a complete copy of the exam.
- The total amount of achievable credits in this simulation is **90 credits**.
- No additional resources are allowed.

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 You are training a network to classify images of handwritten digits in the range of [0,...,9] on the MNIST dataset. Which of the following data augmentation techniques are suitable to use for this task?

- Add Gaussian noise to the images
- Vertically flip the images
- Rotation of the images by 10 degrees
- Change the contrast of the images

✓ 1.2 What is true about Residual Blocks?

- Reduce the number of computations in the forward pass
- Act as a highway for gradient flow
- Enable a more stable training of larger networks
- Act as a regularizer

(↓ ← →)

✓ 1.3 For a fully-convolutional 2D CNN, if we double the spatial dimensions of input images, ...

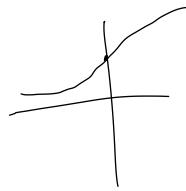
- ... the number of network parameters doubles \times .
- ... the number of network parameters stays the same
- ... the receptive field of an arbitrary pixel in an intermediate activation map can decrease \times
- ... the dropout coefficient p must be corrected to \sqrt{p} in test time \times

✓ 1.4 What is true about Generative Adversarial Networks?

- The Generator minimizes the probability that the Discriminator is correct
- The Generator provides supervision for the Discriminator \times
- The Discriminator acts as a classifier
- The Discriminator samples from a latent space \times .

1.5 Given input x , which of the following statements are always true? Note: For dropout, assume the same set of neurons are chosen.

- $\text{BatchNorm}(\text{ReLU}(x)) \equiv \text{ReLU}(\text{BatchNorm}(x)) \times$
- $\text{Dropout}(\text{ReLU}(x)) \equiv \text{ReLU}(\text{Dropout}(x)) \checkmark$
- $\text{MaxPool}(\text{ReLU}(x)) \equiv \text{ReLU}(\text{MaxPool}(x)) \checkmark$
- $\text{ReLU}(\text{Sigmoid}(x)) \equiv \text{Sigmoid}(\text{ReLU}(x)) \times$
...
...
...
...



1.6 When you are using a deep CNN to train a semantic segmentation model, which of the following can be chosen to help with overfitting issues?

- Decrease the weight decay parameter
- Increase the probability of switching off neurons in dropout \times
- Apply random Gaussian noise to the input images
- Using average pooling layers in the network

1.7 In terms of (full-batch) gradient descent (GD) and (mini-batch) stochastic gradient descent (SGD), which of the following statements are true?

- The computed gradient of the loss w.r.t model parameters in SGD is equal to the computed gradient in GD 使用 mini batch it's
- The expected gradient of the loss w.r.t model parameters in SGD is equal to the expected gradient in GD
- There exists some batch size, for which the gradient of the loss w.r.t model parameters in SGD is equal to the gradient in GD
- SGD and GD will converge to the same model parameters, but SGD requires less memory at the expense of more iterations

1.8 What is true about batch normalization assuming your train and test set are sampled from the same distribution?

- Batch normalization cannot not be used together with dropout doesn't work well
- Batch normalization makes the gradients more stable, so we can train deeper networks
- At test time, Batch normalization uses a mean and variance computed on test set samples to normalize the data \times
- Batch normalization has learnable parameters

1.9 What is true for common architectures like VGG-16 or LeNet? (check all that apply)

- The number of filters tends to increase as we go deeper into the network
- The width and height of the activation maps tends to increase as we go deeper into the network
- The input can be an image of any size as long as its width and height are equal
- They follow the paradigm: Conv \rightarrow Pool ... \rightarrow Conv \rightarrow Pool \rightarrow FC ... \rightarrow FC

Problem 2 Short Questions (18 credits)

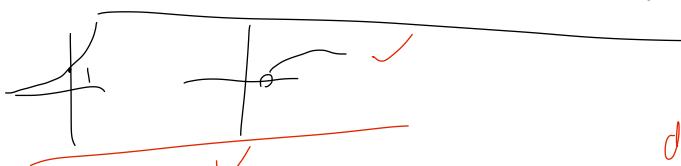
0 2.1 In k -fold cross validation, choosing a larger value for k increases our confidence in the validation score. What could be a practical disadvantage in doing so? Explain how it arises.

- 1 Training time increase, because we need iterate k times for train on k .
 2 Validation set may have different distribution with training set, leading to large $k \uparrow$
- ~~Validation set may have different distribution with training set, leading to large $k \uparrow$~~
- proven model uses data from training set.

0 2.2 Consider the activation function $f: \mathbb{R} \rightarrow \mathbb{R}$ and $f(x) = \ln(1 + e^x)$.

1 Which one of the following activation functions is most closely approximated by f ? Briefly justify your answer (2 points). What is the benefit of f over the activation function it closely approximates (2 points)?

- 2 • Tanh
 3 • ReLU
 4 • Sigmoid



Final decision

ReLU ✓ because $e^x > 0, \forall x$
 $1+e^x > 1, \forall x$
 $\ln(1+e^x) > 0, \forall x$.

can totally represent the real activation function like ReLU.

And even better, being in different regions

only when

$$\text{ReLU}(x) = \max(x, 0)$$

0 2.3 Explain the difference between the validation set and the test set. In your answer, explain the role of each subset and how they are used differently.

1 Step 1 Validation set and test set will all be the subset of data set
 Step 2 Test set is totally new data will generate to in test the or from else, which never seen
 Validation test the generation with different hyperparameters
 generalization
 but test set can be only use once
 When test the final overall part of system
 Validation set will be used any time when test the model and trying the

0 2.4 You notice vanishing/exploding gradients in a deep network using the tanh activation function. Suggest two possible changes you can make to the network in order to diminish this issue, without changing the number of trainable parameters. Explain how each of these changes helps.

- 1 Use ReLU ✓ because ReLU hasn't got gradient problem
 2 Data augmentation
 Residual block
 High-way for gradient
 More data
 drop out, randomly drop the weights, make model less to train

2.5 Can two consecutive dropout layers with probabilities q and p be replaced with one dropout operation? Explain.

0
1
2

✓
Yes

one drop out $(1 - \alpha)$

two drop out work together

$$(1 - (1-p)(1-q)) = p + q - qp$$

✓

2.6 Can one encounter overfitting in an unsupervised learning setting? If your answer is *no*, provide a mathematical reasoning. If your answer is *yes*, provide an example.

0
1
2
3

Yes

Autoencoder has large bottleneck, learn too much information \rightarrow overfitting.

PCA with many components.

2.7 For each of the following functions, describe one common problem when choosing them as the activation function for your deep neural network: (a) Sigmoid, (b) ReLU, (c) Identity

0
1
2
3

(a) Gradient saturation \Rightarrow vanishing gradient

(b) ReLU \Rightarrow not centered / when $x < 0$ $\text{ReLU}(x)=0 \Rightarrow$ dead relu / term never

(c) non linear in the MU

Problem 3 Autoencoder (11 credits)

Consider a given **unlabeled** image dataset consisting of 10 distinct classes of animals.

0 To train an Autoencoder on images, which type of losses you would use? Name two suitable losses.

1
2

Gross Entropy for multiple clas \times L_1, L_2

Binary cross Entropy for ten classes p_{true}

0 To explain the effect of choosing a bottleneck dimension which is too small, and the effect of a too large bottleneck dimension in Autoencoders.

1
2

Too small, AE compresses too much \rightarrow learn too little features or fails to capture input \rightarrow underfitting \rightarrow bad generalization.

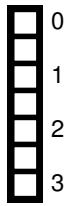
compresses too little \rightarrow learn too much \rightarrow overfitting \rightarrow generalization ~~is same as~~ same as input

0 Having trained an Autoencoder on this dataset, how would you use the trained Autoencoder (without further training/fine-tuning) to partition the dataset into 10 subsets, where each subset consists only of images of a distinct type of animal?

1
2

drop the decoder part and only use encode part
freeze the ~~weights~~ of encoder
add FC layer after the AE, and use softmax func + CE loss
to clarify \times Do clustering

3.4 We want to use the same network architecture for de-noising and colorizing old, degraded, gray-scale images of animals. Given the dataset you already have, explain the steps you would take to train your model. In your answer, elaborate on your model's inputs, outputs, and losses.



Add noisy

Transform the data into gray-scale images

Use L_1, L_2 loss, original image as target

3.5 Explain the differences between Autoencoders and Variational Autoencoders. How do they differ during training?

The bottleneck is code layer,
AE encode the high-dim feature to the low dimension, and
encode sample from the code layer

VAE's encode the high-dim feature to the ~~latent~~ latent space
decode ~~the~~ sample from the lower dim latent space

VAE constraint the latent space distribution into Gaussian distribution

VAE can sample a random vector from latent space to generate images.

AE can't, because without normalize the feature, AE will only generate useless output

AE only need to minimize the loss between original image and generated image
VAE need minimize the MLE.

Problem 4 CNNs (10 credits)

You are given the following network that classifies RGB images into one of 4 classes.

All Conv2d layers use $\text{kernel} = 3$, $\text{padding} = 1$, $\text{stride} = 1$, $\text{bias} = \text{True}$ and are defined as $\text{Conv2d}(<\text{channels}_{\text{in}}>, <\text{channels}_{\text{out}}>)$.

All MaxPool2d layers use $\text{stride} = 2$, $\text{padding} = 0$, and are defined as $\text{MaxPool}(<\text{kernel}>)$.
The input dimension x of the Linear layer is unknown.

The network's architecture is as follows:

- $\text{Conv2d}(3, 8) \rightarrow \text{MaxPool2d}(2) \rightarrow \text{BatchNorm2d()} \rightarrow \text{ReLU()} \rightarrow$
- $\text{Conv2d}(8, 16) \rightarrow \text{MaxPool2d}(2) \rightarrow \text{BatchNorm2d()} \rightarrow \text{ReLU()} \rightarrow$
- $\text{Conv2d}(16, 32) \rightarrow \text{MaxPool2d}(2) \rightarrow \text{BatchNorm2d()} \rightarrow \text{ReLU()} \rightarrow$
- $\text{Flatten()} \rightarrow$
- $\text{Linear}(x, 4) \rightarrow \text{Softmax()}$

0 4.1 In terms of x , what is the total number of trainable parameters of the last linear layer? Include a bias term in your calculation.

1
2

$$4x + 4$$

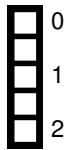
0 4.2 Given RGB input images of size 80×80 pixels, what should the value of x in the Linear layer be?
Explain your calculation.

1
2

$$\begin{aligned} \frac{F_{\text{in}} + 2P - N}{S} + 1 &= F_{\text{out}} & \frac{80 + 2 \times 0 - 3}{1} + 1 &= 80 \times 80 \times 8 & \frac{40 + 2 \times 0 - 3}{2} + 1 &= 20 \times 20 \times 16 \\ \frac{80 + 2 \times 0 - 2}{2} + 1 &= 40 \times 40 \times 8 & \frac{20 + 2 \times 0 - 3}{1} + 1 &= 20 \times 20 \times 32 \\ \frac{40 + 2 \times 0 - 3}{1} + 1 &= 40 \times 40 \times 16 & \frac{20 + 2 \times 0 - 3}{1} + 1 &= 10 \times 10 \times 32, \\ x &= 10 \times 10 \times 32 & & & \approx 3200 \end{aligned}$$

4.3 Explain the main difference between the usage of a BatchNorm layer in a convolutional network in comparison to a fully connected network.

BN in FC : Normalize the each input. on Channel dimension,
BN in CNN : Normalize the each feature map



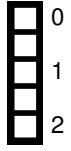
4.4 Compute the total number of trainable parameters of the first convolutional layer, Conv2d(3,8).

$$\text{Weight Tensor } (3, 8, 3, 3) \quad (3 \times 3 \times 3 + 1) \times 8 = 224$$



4.5 Compute the total number of trainable parameters in all of the BatchNorm layers.

$$8 \times 2 + (b \times 2 + 32 \times 2) = 112$$



Problem 5 Optimization and Gradients (16 credits)

You are training a large fully-connected neural network and select as an initial choice an SGD optimizer. In order to overcome the limitations of SGD, your colleague suggests adding momentum.

- 0 5.1 Name two limitations of SGD that momentum can potentially solve. Explain how momentum solves them.

Avoid stuck at saddle point / speeds up the convergence / solve some problems
Momentum introduce a exponentially weight moving average over the gradient
slow / small step) can escape local minimum / SGD is noisy / only one family members
Speed up learning if gradient toward to same direction /
keep direction keep escape local minimum

- 0 5.2 One can apply momentum, as shown in the formula:

$$\nu^{k+1} = \beta \cdot \nu^k - \alpha \cdot \nabla_{\theta} L(\theta^k)$$

- 2 What do the hyperparameters α and β represent?

β accumulation rate of velocity
 α learning rate

- 0 5.3 How does Nesterov Momentum differ from standard momentum? Explain.

It called Look-ahead - method, which is a second momentum
allow large learning rate permit the next gradient step and jump to it
Gradient computed from previous calculated gradient and do correction.
Gradient correction / avoid choosing learning rate
in case over shoot

- 0 5.4 Is RMSProp considered a first or second order method (1p)? What is the main difference between RMSProp and SGD+Momentum?

Second Moment first with learning rate divide exponentially decay average of square gradient
Allow large learning rate / avoid chose learning rate

For the following questions, consider the convex optimization objective:

$$\min_{x \in \mathbb{R}} x^2$$

5.5 What is the optimal solution of this optimization problem?

F.O.C $\frac{d x^2}{dx} = 2x = 0$
 $x=0$

0
1

5.6 You are working with an initialization of $x_0 = 5$ and a learning rate of $lr = 1$. How many iterations would gradient descent (without momentum) need in order to converge to the optimal solution? Explain.

$$\begin{aligned} x^+ &= x - lr \cdot \nabla(x) \\ x^+ &= 5 - 1 \cdot (2 \cdot 5) \\ &= -5 \end{aligned}$$

$x^+ = -5 - 1 \cdot (-10)$
 $= 5$
Never converge ✓

0
1

5.7 Assuming you instead start with a random initialization of x_0 , how could you speed up the convergence of the gradient descent optimizer (without adding momentum) in this case?

Use random Initialization

$$\text{Var}(n) = \frac{1}{n}$$

Reduce adaptive learning rate
dynamic

0
1

5.8 What is the main advantage of using a second order method such as Newton's Method? Why are second order methods not used often in practice for training deep neural networks?

Fast convergence \Leftrightarrow less iterations

No need for slow learning rate

Compute Inverse Hessian matrix is computationally cost

Second order doesn't work well in mini batch

0
1
2

5.9 How many iterations would Newton's method need to converge (using the same initialization $x_0 = 5$, $lr = 1$)? Explain.

$$x^+ = x -$$

0
1

Problem 6 Derivatives (9 credits)

Consider the formula of the Sigmoid function $\sigma(x) : \mathbb{R} \rightarrow \mathbb{R}$:

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad (\cancel{1+e^{-x}})^{-1}$$

0 6.1 Compute the derivative $\frac{d\sigma(x)}{dx}$ in terms of x .

$$\begin{aligned}\frac{d\sigma(x)}{dx} &= (-1) \cdot (1 + e^{-x})^{-2} \cdot (e^{-x}) \cdot (-1) \\ &= \frac{e^{-x}}{(1 + e^{-x})^2}\end{aligned}$$

0 6.2 A special property of this function is that its derivative can be expressed in terms of the Sigmoid function itself. Denote $y = \sigma(x)$, and show how the derivative you computed can be re-written in terms of y , the output of the Sigmoid function. Hint: Your answer should only depend on y .

$$\frac{dy}{dx} = y(1-y) \quad 1-y = \frac{e^{-x}}{1+e^{-x}} \quad y = \frac{1}{1+e^{-x}}$$

An affine Layer is described by $\mathbf{z} = \mathbf{XW} + \mathbf{b}$.

Consider the following affine layer, which has 2 input neurons and 1 output neuron:

$$\mathbf{W} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}_{2 \times 1} \quad \begin{matrix} \textcircled{1} \\ \textcircled{2} \\ \textcircled{3} \end{matrix}$$

$$\mathbf{b} = 2 \in \mathbb{R}^1$$

and input:

$$\mathbf{X} = \begin{bmatrix} 1 & 1 \\ 0 & -1 \end{bmatrix}_{2 \times 2}$$

The forward pass of the network would be:

$$\sigma(\mathbf{z}) = \sigma(\mathbf{XW} + \mathbf{b}) = \sigma\left(\begin{bmatrix} 1 & 1 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} + 2\right) = \sigma\left(\begin{bmatrix} 3 \\ -2 \end{bmatrix} + \begin{bmatrix} 2 \\ 2 \end{bmatrix}\right) = \sigma\left(\begin{bmatrix} 5 \\ 0 \end{bmatrix}\right) = \begin{bmatrix} 1 \\ 0.5 \end{bmatrix} \text{ (rounded up).}$$

Let's compute the backward pass of the network.

6.3 If $\mathbf{y} = \sigma(\mathbf{z}) = \begin{bmatrix} 1 \\ 0.5 \end{bmatrix}$, calculate the gradient of the output after the Sigmoid activation function

w.r.t \mathbf{z} , $\frac{dy}{dz}$:

$$\frac{dy}{dz} = y(1-y) = \begin{bmatrix} 1 \\ 0.5 \end{bmatrix} \circ \begin{bmatrix} 0 \\ 0.5 \end{bmatrix} = \begin{bmatrix} 0 \\ 0.25 \end{bmatrix}$$

6.4 We will use the computed gradient to perform back-propagation through the affine layer to the network's parameters.

Let $dout$ be the upstream derivative of the Sigmoid that you have calculated in question 3. Calculate the derivatives $\frac{dy}{dw}$ and $\frac{dy}{db}$.

Hint: Pay attention to the shapes of the results; they should be compatible for a gradient update.

Note: In case you skipped the previous question, you can get partial points by writing the correct formulas using $dout$ symbolically.

$$\frac{dy}{dw} = \frac{dy}{dz} \cdot \frac{dz}{dw} = \begin{bmatrix} 0 \\ 0.25 \end{bmatrix} \circ \begin{bmatrix} 1 & 1 \\ 0 & -1 \end{bmatrix} = \begin{bmatrix} 0 \\ -0.25 \end{bmatrix}$$

$$\frac{dy}{db} = \frac{dy}{dz} \cdot \frac{dz}{db} = \begin{bmatrix} 0 \\ 0.25 \end{bmatrix} \circ \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0.25 \end{bmatrix}$$

$$\frac{dy}{dw} = \mathbf{X}^T \cdot dout = \begin{bmatrix} 1 & 1 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} 0 \\ 0.25 \end{bmatrix} = \begin{bmatrix} 0 \\ -0.25 \end{bmatrix}$$

$$\frac{dy}{db} = b \cdot dout = \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 0.25 \end{bmatrix} = 0.25$$

0
1
2
3

0
1
2
3
4

Problem 7 Model Evaluation (8 credits)

Two students, *Erika* and *Max* train a neural network for the task of image classification. They use a dataset which is divided into train and validation sets. They each train their own network for 25 epochs.

- 0 7.1 Erika selects a model and obtains the following curves. Interpret the model's behaviour from the curves. Then, suggest what could Erika do in order to improve its performance?

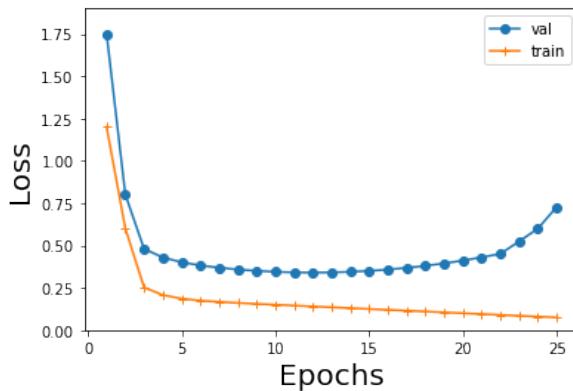


Figure 7.1: Training curves for Erika's model.

After 10 epochs , val loss can increase , That means Overfitting
weight decay regularization
Dropout
Not enough
early stopping

- 0 1 2 7.2 Max selects a different model and obtains the following curves. Interpret the model's behaviour from the curves. Then, suggest what change could Max make to his model in order to improve its performance?

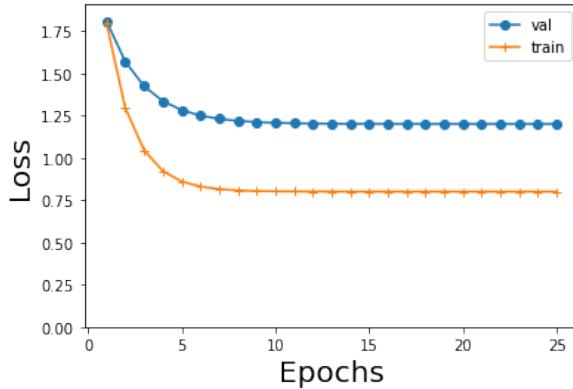


Figure 7.2: Training curves for Max's model.

Two losses are very high , reduce the learning rate
Underfitting
by Dropout . BN

7.3 Both Max and Erika are able to agree on a model architecture and obtain the following curves. However, when deployed in real world, their model seems to perform poorly. What is a possible reason for such an observation and what should they do?

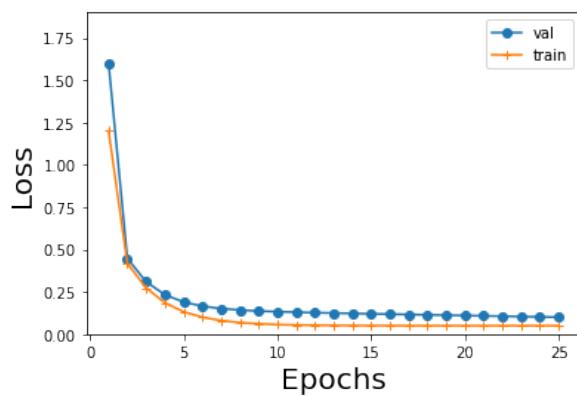


Figure 7.3: Training curves for the new model.

Test set and Train set may have different distribution
dis match

~~too much weight~~, model learn something from train set, can not generalize good

Data converges more slowly

After adapting the new network architecture, Max and Erika are training their own model, using the same architecture, with identical initial weights, using exactly the same hyperparameters. They also use the same SGD optimizer (no momentum), batch size, and learning rates. The only difference is that Max normalizes the loss by $1/N$ (where N is the number of training samples in the dataset) while Erika does not.

7.4 How does this affect the optimal model weights that minimize this optimization objective? (1p) After 10 optimizer steps, will they arrive at the same model parameters? Explain.(2p)

the gradient is more stable
because BN normalizes the input data

BN doesn't affect the open weights

IVD search (GD). Stochastic

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.