



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. Dr. Matthias Nießner

Date: Wednesday 19th February, 2020
Time: 13:30 – 15:00

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

Working instructions

- This exam consists of **20 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 exam is 89 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.

Left room from _____ to _____ / Early submission at _____

Problem 1 Multiple Choice (18 credits)

Multiple Choice Questions:

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

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 put an additional **cross** besides the box.

~~check~~

a) You train a neural network and the loss diverges. What are reasonable things to do?

- Try a different optimizer. 
 - Add dropout.
 - Decrease the learning rate.
 - Increase the number of parameters.
- Check for proper initialization of weights.
Lower the learning rate.
Check for outliers in the data that may be causing the model to diverge.
Check for missing or corrupted data in the training set.
Try a different optimization algorithm.
Make sure the input data is scaled properly.
Increase the batch size.
Use regularization techniques such as dropout or weight decay to prevent overfitting.
Increase the number of hidden layers or units in the model.
Try a different architecture such as a Convolutional Neural Network (CNN) or Recurrent Neural Network (RNN).

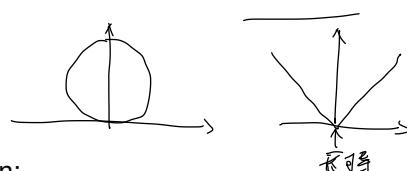
b) Use a neuron to regress the **AND function**. The activation function of the neuron is $f(x) = 0$ if $x < 0$ else 1. What is the weight and bias of the neuron?

x_1	x_2	$x_1 \text{ AND } x_2$	1.0	0.5	0.0	-0.5	-1.0
1.0	1.0	1.0	1.0	0.5	0.0	-0.5	-1.0
1.0	0.0	0.0	1.0	0.5	0.0	-0.5	-1.0
0.0	0.0	0.0	1.0	0.5	0.0	-0.5	-1.0
0.0	1.0	0.0	1.0	0.5	0.0	-0.5	-1.0

- Bias: 1.5, w_1 : 2.0, w_2 : 2.0
- Bias: -1.0, w_1 : 1.5, w_2 : 1.5
- Bias: -1.5, w_1 : 1.0, w_2 : 1.0
- Bias: -2.0, w_1 : 1.4, w_2 : 1.2

c) You want to train an autoencoder to overfit a single image with a fixed learning rate. Setting numerical precision aside, which loss function is able to reach zero loss after training with gradient descent?

- L2
- L1



d) Regularization:

- Is a technique that aims to reduce the generalization gap.
- Dropout, the use of ReLU activation functions, and early stopping can all be considered regularization techniques.
- Weight decay (L^2) is commonly applied in neural networks to spread the decision power among as many neurons as possible.
- Is a technique that aims to reduce your validation error and increases your training accuracy.

e) Which statements are correct for a Rectified Linear Unit (ReLU) applied in a CNN?

- Despite a small learning rate, without saturation the gradients are very likely to explode.
- A large negative bias in the previous layer can cause the ReLU to always output zero.
- Large and consistent gradients allow for a fast network convergence.
- Max pooling must always be applied after the ReLU.

f) You want to use a convolutional layer to decrease your RGB image size from 254x254 to 127x127. What parameter triplets achieve this?

- Kernel size 3, stride 2, padding 1
- Kernel size 6, stride 2, padding 2
- Kernel size 2, stride 2, padding 0
- Kernel size 5, stride 2, padding 2

$$\begin{array}{ccccc}
 & 4 & & F + 2p - n & +1 \\
 & 6 & & S & \\
 & 2 & & 254 + 2p - n = 127 \times 2 = & 252 \\
 & 6 & & 2p - n = 2 & \\
 & & & n = 2p + 2 &
 \end{array}$$

g) Your train and val loss converge to about the same value. What would you do to increase the performance of your model?

Add more input features.

Increase the capacity of your model.

Collect more data

Add more training data. ?

Increase model complexity

Add regularization.

Feature engineering: PCA

Hypertuning __'

h) An autoencoder

has no loss function.

has a bottle neck layer.

指 transfer learning

is often used for fine-tuning pre-trained models.

用已训练模型 → model

requires class labels for training.

i) 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

e d b a c

ebadc

eadbc

edbac

Problem 2 Short Questions (22 credits)

ReLU kill half of data

- 0 a) Kaiming initialization corresponds to Xavier initialization with the variance multiplied by two. In which case (1p) and why (1p) would you chose this initialization?

1 When we use ReLU as activation / because if we use Xavier and ReLU
2 there still have the same problem
Gradient vanishing

- 0 b) You are given a convolutional layer with kernel size 3, number of filters 3, stride 1 and padding 1. Compute the shape of the weights (0.5p) and write them down explicitly such that this convolutional layer represents the identity for an RGB image input (1.5p).

1 (3 3 3 3)

?

- 0 c) Name two reasons to use an inception layer in favor of a standard convolutional layer.

1 Don't know which filter to use
2 Spatial Size reduction
Compute efficient

- 0 d) What are the definitions of *bias* and *variance* in the context of machine learning?

1 The difference of the data
2 Bias : Underfitting / Difference between Ground truth error and training set error
Variance : Overfitting / Difference between training set error and validation/test set error.

e) In Generative Adversarial Networks the generator and discriminator play a two player min-max game. Why is this hard to optimize and what heuristical method is used instead of the default min-max formulation.

G minimizes the probability when D is correct G can not learn when D reject all generated sample

$D(G(z))$

Want 0 want 1

Equilibrium is saddle point

“Alternating optimization” Fixed Generator update via versa.

f) Consider the quote below. Demonstrate how a fully connected layer with an input size of 512 neurons and an output of 10 neurons can be modeled as a convolutional layer.



Yann LeCun

6. April 2015 ·

$512 \times 1 \times 1 \rightarrow 10 \times 1 \times 1$
1x1 convolution layer

In Convolutional Nets, there is no such thing as "fully-connected layers".
There are only convolution layers

use 10 filter with input channel 512
 $512 \rightarrow 1 \times 1 \times 512 \xrightarrow{\text{use } 1 \times 1 \text{ convolution layer}} 1 \times 1 \times 10$
 Weights $E(512, 10, 1, 1)$

g) Explain the Markov Assumption in reinforcement learning.

The future is independent of the past given the present

$$P(S_{t+1} | S_t) = P(S_{t+1} | S_t, S_{t-1}, \dots, S_1)$$

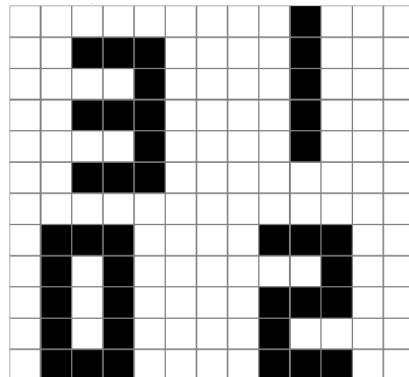
h) Where do we use neural networks in Q-learning (1p) and why are they needed (1p)?

Use Q-learning to approximate Q-value

To estimate how good the state is

- 0 
 i) The following image shows a rectangular image of size 16×16 . Design a 5×5 convolutional filter that is maximally activated when sliding over the '3' in the image below (Black pixels are 1s and white are -1. For simplicity use only -1s and 1s in your designed filter).

1
 2



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

- 0 
 j) Name one advantage and one disadvantage of Recurrent Neural Networks in general.

1
 2

Deal with the sequential data

Without non-linearity

Vanishing gradient problem

- 0 
 k) Long-Short Term Memory Units suffer less from the vanishing gradient problem than vanilla RNNs. What two design changes make this possible?

1
 2

Add cell to transport the information.

Add three gates and use tanh as non-linearity
gate system

Problem 3 House Prices and Backpropagation (11 credits)

You are tasked to predict house prices for your first job, i.e. for a given input vector of numbers you have to predict a single floating point number that indicates the house price (e.g., between 0 and 1m euros).

a) What network loss function would you suggest for this problem (0.5p) and why (0.5p)?

Cross entropy

because it's a multi-classification

MSE

Floating point regression

	0
	1

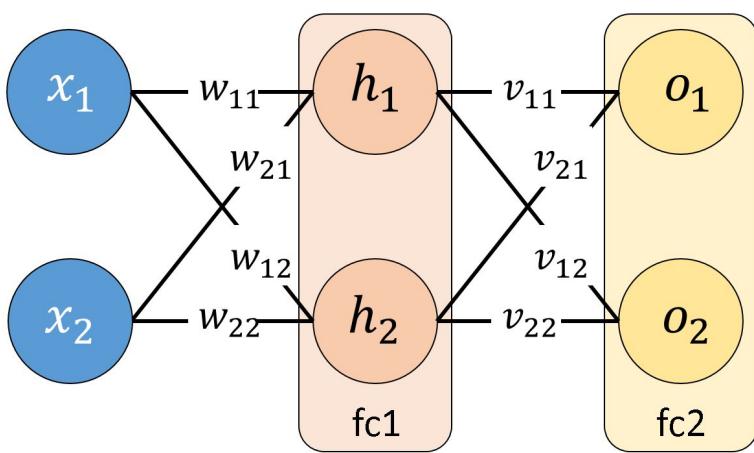
b) How would you approximate the task as a classification problem (0.5) and which loss function would you propose in this situation (0.5)?

Bucket price into price range

Cross Entropy

	0
	1

After collecting some data you start off with a simple architecture



Your architecture is composed of two fully-connected layers (fc1 and fc2) which both contain

- a linear layer with weights and biases as outlined on the next page,
- followed by a Dropout with a probability parameter of 0.5,
- as well as a Leaky ReLu with a parameter of 0.5 as your non-linearity of choice at the end.

c) How many trainable parameters does your network have?

$$2 \times 2 + 2 + 2 \times 2 + 2 = 12$$

	0
	1

Weights and biases of the linear layers:

Layer	fc1			fc2		
Nodes	h_1		h_2	o_1		o_2
Variable	w_{11}	w_{21}	b_{h_1}	w_{12}	w_{22}	b_{h_2}
Value	1.0	-1.0	1.0	-1.0	1.0	-1.0

- 0 d) You are experiencing difficulties during training and thus decide to check the network in test mode manually.
In your test case the input x values are

$$1 \quad x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}. \\ 2$$

3 What is the output of your network for that input? Write down all computations for each step.

$$4 h_1 = w_{11} \cdot x_1 + w_{21} \cdot x_2 + b_{h_1} = 1 \cdot 1 + (-1) \cdot 0 + 1.0 = 2 \quad \text{drop out} = 1$$

$$h_2 = w_{12} \cdot x_1 + w_{22} \cdot x_2 + b_{h_2} = (-1.0) \cdot 1 + (1.0) \cdot 0 + (-1.0) = -2 \quad = -1$$

$$\text{leaky relu } (h_1) = 1 = z_1$$

$$\text{linear relu } (h_2) = -0.5 = z_2$$

$$O_1 = v_{11} \cdot z_1 + v_{21} \cdot z_2 + b_{o_1} = 0.7 \cdot 1 + (-1.0) \cdot (-0.5) + 1.0 = 2 \rightarrow 1$$

$$O_2 = v_{12} \cdot z_1 + v_{22} \cdot z_2 + b_{o_2} = 1.0 \cdot 1 + (1.0) \cdot (-0.5) + 0.0 = 0.5 \rightarrow 0.25$$

$$\text{linear } (o_1) = 1 = z_3$$

$$\text{linear } (o_2) = 0.25 = z_4$$

$$z_o = \begin{pmatrix} 1 \\ 0.25 \end{pmatrix}$$

- 0 e) As you were unsure of your loss function, you phrase the task as a binary classification problem for each of your two outputs independently. Calculate the binary cross-entropy with respect to the natural logarithm
1 for the labels

$$y = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}.$$

You may assume $0 \cdot \ln(0) = 0$. (Write down the equations and keep the simplified logarithm.)

$$\text{BCE : } - (y_1 \ln(\hat{y}_1) + (1-y_2) \ln(1-\hat{y}_2))$$

$$= - (1 \cdot \ln(1) + 1 \cdot \ln(0.75))$$

$$= - \ln(0.75)$$

f) Please update the weight v_{21} using gradient descent with a learning rate of 1.0 with respect to the binary cross-entropy loss as well as labels from e) and the forward computation from d). (Please write down all your computations. Writing down formulas and values in general form is encouraged)

$$\begin{aligned}
 \frac{\partial BCE}{\partial v_{21}} &= \underbrace{\partial - (y_1 \ln(z_3) + (1-y_2) \ln(z_4))}_{\partial z} \cdot \underbrace{\frac{\partial z_3}{\partial dropout}}_{\partial o_1} \cdot \underbrace{\frac{\partial dropout}{\partial o_1}}_{\partial v_{21}} \cdot \underbrace{\frac{\partial o_1}{\partial v_{21}}} \\
 &= - \frac{y_1}{z_3} \cdot 1 \cdot 0.5 \cdot z_2 \\
 &= - | . | \cdot 0.5 \cdot \sim 0.5 \\
 &= 0.25 \\
 v_{21}^+ &= v_{21} - \alpha \cdot \frac{\partial BCE}{\partial v_{21}} \\
 &= 1.0 - 1.0 \cdot 0.25 \\
 &= 1.25
 \end{aligned}$$

	0
	1
	2
	3

Problem 4 Optimization (11 credits)

0 a) Why can't we expect to find a global minima while training neural networks?

1 because usually it's not a convex problem case there exist some local minima
many local minima
no practical way to say which is global minima

0 b) Why is finding a local minimum often enough? ^{many} Local minima provide some performance

1 In some cases, when it's a convex problem ^{many} Local minima is global minima
when generalization

In other situations it's hard to find a global minimum

0 c) What is a saddle point (1p)? What is the advantage/disadvantage of Stochastic Gradient Descent (SGD) in dealing with saddle points (1p)?

1 A saddle point is the in this point, gradient is 0
but this point is neither local min nor local ~~max~~ max

Because of the noisy updates and "stochastic" SGD may have
chance to leave saddle point and continue ~~not~~ find optimum

Disadvantage - no momentum

0 d) Explain the concept behind momentum in SGD.

1 Momentum : Introduce a technique which can adapt the learning rate,
that can solve the scale problem in SGD
~~a void stuck in local minima~~ speedup optimization

0 e) Which optimizer introduced in the lecture uses second but not first order momentum?

1 ~~Look-ahead-Momentum~~ No RMSProp

0 f) Explain the beta (β_1 and β_2) hyperparameters of Adam with respect to the gradients.

1 Two exponentially decay parameters
 β_1 is for first-order Momentum estimate
 β_2 is for second-order Momentum estimate

g) What would be an advantage of a second order optimization method such as the Newton method (1p) ? Why is it not commonly used in the context of neural networks?

0
1

Faster than First-order Momentum

It's computational costly for H^{-1}

Stochasticity work doesn't well with second order them

h) Name the key idea of Newton's method (1p) and write down the update step formula (1p).

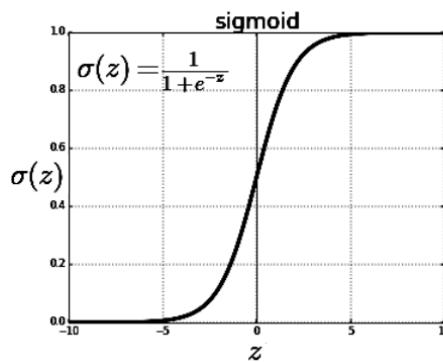
0
1
2

Use Hessian matrix (second order derivative) which can speed up the converge process Approximate second order Taylor expansion

$$\theta' = \theta - H^{-1} \nabla L$$

Problem 5 Network Architectures and Training (10 credits)

You are training a neural network with 10 convolutional layers and the activation function shown below:



- 0 a) What is the purpose of activation functions in neural networks?

1 Add non-linearity to NN such that NN can deal with non-linear data

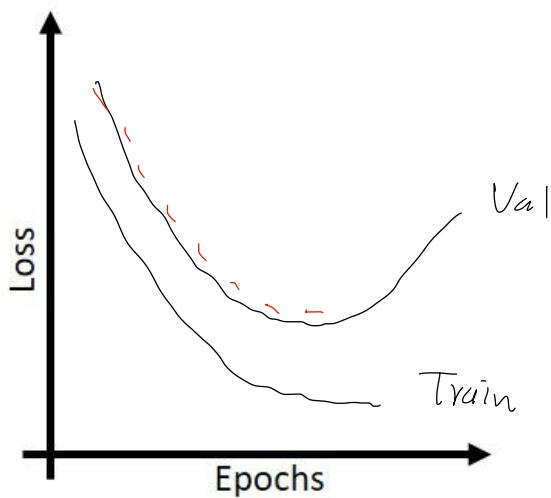
- 0 b) You find that your network still is not training properly and discover that your network weights have all been default initialized to 1. Why might this cause issues for training?

If the layer have same weight, then they get same gradient via backprop
That means NN learn the same thing every time

- 0 c) How can you resolve the problems with weight initialization and provide stable training in most scenarios for the proposed network? Explain your solution.

1 Use Xavier Initialization
which mean is 0 and variance is $\frac{1}{n}$
ensure the variance of input and output are same

- 0 d) After employing your solutions, you are ready to train your network for image segmentation. After 50 epochs, you come to the conclusion that the network is too large for such a task. What is this effect called? How do you make this observation? Make a plot of the corresponding training (regular line) and validation losses (dashed line) and name them appropriately.



Overfitted ✓

Use graph to show the train and validate loss
check the difference between train and Val

- e) Without changing the convolutional layers of your network, name two approaches to counteract the problems encountered in d).

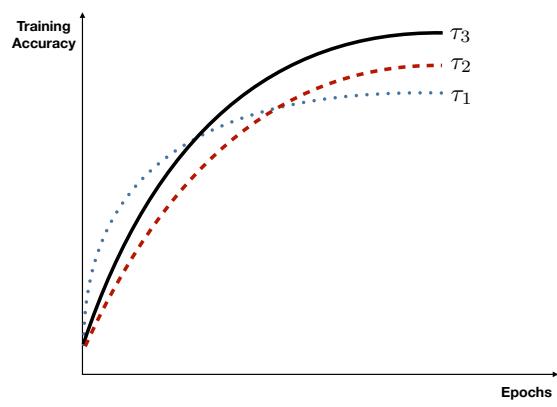
Use Adam as optimizer to adapt the learning rate

Regularization / dropout / data augmentation / more data
(early stopping)

0
1
2

- f) You adapt your network training accordingly, and now are performing a grid search to find the optimal hyperparameters for vanilla stochastic gradient descent (SGD). You try three learning rates τ_i with $i \in \{1, 2, 3\}$, and obtain the following three curves for the training accuracy. Order the learning rates from larger to smaller.

0
1
2



$$\tau_1 > \tau_3 > \tau_2$$

Problem 6 Batchnormalization (5 credits)

A friend suggested you to use Batchnormalization in your network. Recall the batch normalization layer takes values $x = (x(1), \dots, x(m))$ as input and computes $x_{norm} = (x_{norm}^{(1)}, \dots, x_{norm}^{(m)})$ according to:

$$x_{norm}^{(k)} = \frac{x^{(k)} - \mu}{\sqrt{\sigma^2}} \quad \text{where } \mu = \frac{1}{m} \sum_{k=1}^m x^{(k)}, \quad \sigma^2 = \frac{1}{m} (x^{(k)} - \mu)^2.$$

It then applies a second transformation to get $y = (y^{(1)}, \dots, y^{(m)})$ using learned parameters $\gamma^{(k)}$ and $\beta^{(k)}$:

$$y^{(k)} = \gamma^{(k)}(x_{norm}^{(k)}) + \beta^{(k)}.$$

a) How would you make the formulation above numerically stable?

Add a noise to the variance

$$\tilde{x}^{(k)} = \frac{x^{(k)} - E(x^{(k)})}{\sqrt{\text{Var}(x^{(k)}) + \epsilon}} \quad \text{Incase Variance is 0}$$

0
1

b) How can the network undo the normalization operation of Batchnorm? Write down the exact parameters.

Changing two parameters γ and β

$$\gamma^{(k)} = \sqrt{\text{Var}(\tilde{x}^{(k)})} \quad \beta^{(k)} = E(\tilde{x}^{(k)})$$

0
1
2

c) Name the main difference of Batchnormalization during training or testing and note down eventual parameters that need to be stored.

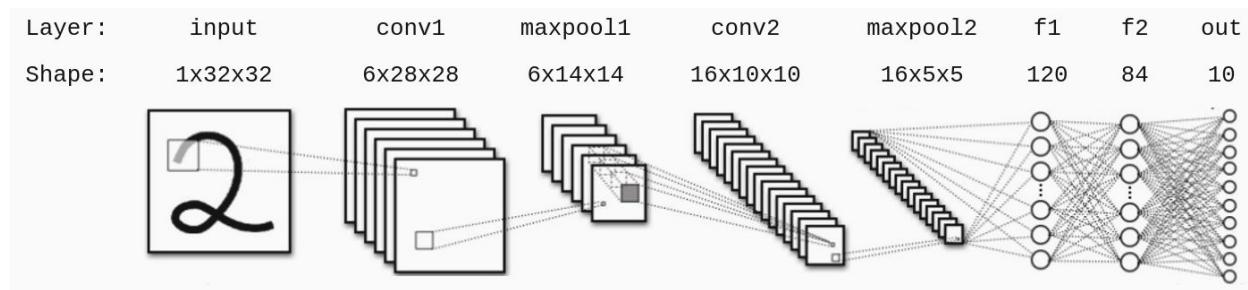
Train Compute mean and variance from mini-batch, and store them

0
1
2

Test Use stored mean and variance

Problem 7 Convolutional Neural Networks (12 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 6x28x28.

All operations are done with stride 1 and no padding. For the convolution layers, assume a kernel size of 5x5.

- 0 a) Explain the term ‘receptive field’ (1p). What is the receptive field of one pixel after the operation ‘maxpool1’ (1p)? What is the receptive field of a neuron in layer ‘f1’ (1p)? *affected by*

1 After convolutional operations, a output pixel can observe how many input pixels

2 For pooling $w = \frac{W-F}{S} + 1$ $14 = \frac{28-2}{2} + 1$ $1 = \frac{W_{in}-2}{2} + 1$ $W_{in} = 2$

3 Conv $w = \frac{N+2P-F}{S} + 1$ $2 = \frac{N+0-5}{1} + 1$ $N = 6$ ~~$b \times b$~~

Whole image 32x32 ~~$16 \times 5 \times 5 = 120$~~

- 0 b) Instead of digits, you now want to be able to classify handwritten alphabetic characters (26 characters). What is the **minimal** change in network architecture needed in order to support this?

1 Last FC, which need to adapt to 26

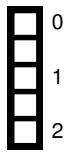
- 0 c) Instead of taking 32×32 images, you now want to train the network to classify images of size 68×68 . List two possible architecture changes to support this?

1 Add a conv layer to adapt the input size *Preprocess \leftrightarrow downsample*
conv 5x5 + maxpooling 2x2

2 or change the first layer
Change the F1 input layer

Fully convolution Layer + global average pooling

d) Your architecture works and you manage to classify characters fairly well. After reading many online blogs, you decide to try out a much deeper network to boost the network's capacity. Name 2 problems that you might encounter when training very deep networks.



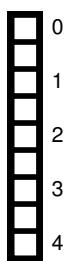
Overfitting, NN learn all the failure from the training set, but can't perform good generalization on unseen data

Memory issues

Gradient vanishing, the deeper the NN is, the gradient is more likely to become closer to zero and update the parameters very slow.

Increasing training time

e) You read that skip connections are beneficial for training deep networks. In the following image you can see a segment of a very deep architecture that uses skip connections. How are skip connections helpful? (1p). Demonstrate this mathematically by computing gradient of output z with respect to w_0 for the network below in comparison to the case without a skip connection (3p). For simplicity, you can assume that gradient of ReLU, $\frac{d(\text{ReLU}(p))}{dp} = 1$.



prevent gradient vanishing / high way for backward

$$z = \text{ReLU}(z')$$

$$z' = G(y) + y$$

$$G(y) = w_3 \cdot \text{ReLU}(w_2 \cdot y)$$

$$y = \text{ReLU}(y')$$

$$y' = F(x) + x$$

$$F(x) = w_1 \cdot \text{ReLU}(w_0 \cdot x)$$

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

$$= 1 \cdot \left(\frac{dG(y)}{dy} + 1 \right) \cdot 1 \cdot \frac{dF(x)}{dw_0}$$

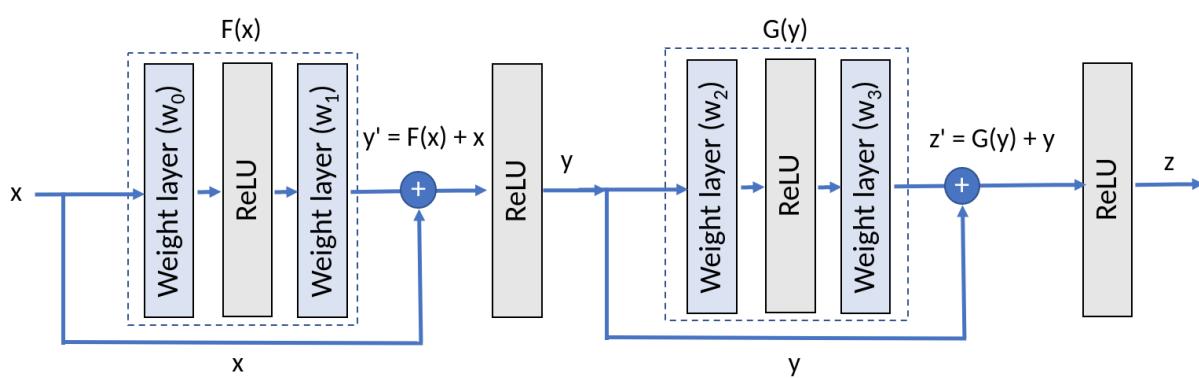
$$= (w_3 \cdot 1 \cdot w_2 + 1) \cdot (w_1 \cdot 1 \cdot x)$$

$$= w_1 \cdot w_2 \cdot w_3 \cdot x + w_1 \cdot x$$

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

$$= 1 \cdot \left(\frac{dG(y)}{dy} \right) \cdot 1 \cdot \frac{dF(x)}{dw_0}$$

$$= (w_3 \cdot 1 \cdot w_2) \cdot (w_1 \cdot 1 \cdot x) = w_1 \cdot w_2 \cdot w_3 \cdot x$$



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.

