

Problem 1 Multiple Choice (18 credits)

Below you can see how you can answer multiple choice questions.

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



- 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. Which of the following models are unsupervised learning methods?

- Auto-Encoder
- Maximum Likelihood Estimate
- K-means Clustering
- Linear regression

1.2 In which cases would you usually reduce the learning rate when training a neural network?

- When the training loss stops decreasing
- To reduce memory consumption
- After increasing the mini-batch size
- After reducing the mini-batch size

1.3 Which techniques will typically decrease your **training** loss?

- Add additional training data
- Remove data augmentation
- Add batch normalization
- Add dropout

1.4 Which techniques will typically decrease your **validation** loss?

- Add dropout
- Add additional training data
- Remove data augmentation
- Use ReLU activations instead of LeakyReLU

DO NOT SCAN/UPLOAD

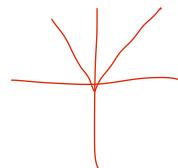
$$\nabla \hat{L} = \nabla L - \alpha \cdot \nabla \ell$$

1.5 Which of the following are affected by multiplying the loss function by a constant positive value when using SGD?

- Memory consumption during training
- Magnitude of the gradient step
- Location of minima
- Number of mini-batches per epoch

1.6 Which of the following functions are not suitable as activation functions to add non-linearity to a network?

- $\sin(x)$
 - $\text{ReLU}(x) - \text{ReLU}(-x)$
 - $\log(\text{ReLU}(x) + 1)$
 - $\log(\text{ReLU}(x + 1))$
- ReLU(0)*
- $x > 0 = x$
 $x = 0$
 $x < 0 = -x$
- $\Rightarrow \ln(0)$ *无定义*



1.7 Which of the following introduce non-linearity in the neural network?

- LeakyReLU with $\alpha = 0$
- Convolution
- Batch Norm
- Skip connection

1.8 Compared to the L1 loss, the L2 loss...

- is robust to outliers *L1 Loss 更 robust*
- is costly to compute *L1 计算绝对值更费力*
- has a different optimum *不同优化*
- will lead to sparser solutions *L1 → ↗*

1.9 Which of the following datasets are NOT i.i.d. (independent and identically distributed)?

- A sequence (toss number, result) of 10,000 coin flips using biased coins with $p(\text{toss result} = 1) = 0.7$
- A set of (image, label) pairs where each image is a frame in a video and each label indicates whether that frame contains humans.
- A monthly sample of Munich's population over the past 100 years
- A set of (image, number) pairs where each image is a chest X-ray of a different human and each number represents the volume of their lungs.

DO NOT SCAN/UPLOAD

Problem 2 Short Questions (29 credits)

0
1
2
3

- 2.1 Explain the idea of data augmentation (1p). Specify 4 different data augmentation techniques you can apply on a dataset of RGB images (2p).

DO NOT SCAN/UPLOAD

0
1
2

- 2.2 You are training a deep neural network for the task of binary classification using the Binary Cross Entropy loss. What is the expected loss value for the first mini-batch with batch size $N = 64$ for an untrained, randomly initialized network? Hint: $BCE = -\frac{1}{N} \sum y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i)$

DO NOT SCAN/UPLOAD

0
1
2

- 2.3 Explain the differences between *ReLU*, *LeakyReLU* and *Parametric ReLU*.

DO NOT SCAN/UPLOAD

0
1
2

- 2.4 How will weights be initialized by Xavier initialization? Which mean and variance will the weights have? Which mean and variance will the output data have?

DO NOT SCAN/UPLOAD

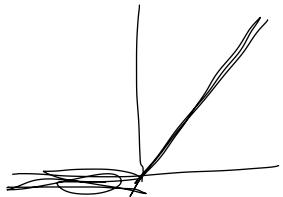
DO NOT SCAN/UPLOAD

2.1 Use different method to process the train data and add this processed data to the origin data

Filping / Cropping / Gaussian blur / Sharpen / Rotation /

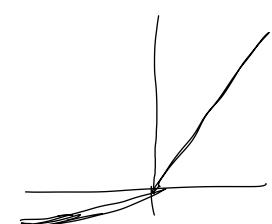
2.2

2.3 ReLU



$$\max(0, x)$$

Leaky ReLU



$$\max(0.1x, x)$$

Piecewise ReLU

$$\max(\alpha x, x)$$

2.4 $\text{Var}(w) = \frac{1}{n}$ $\text{Mean} = 0$

Xavier Initialization ensure the variance of outputs are same as input

same as input

2.5 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.

0
1
2
3

2.6 Given a Convolution Layer in a network with 6 filters, kernel size 5, a stride of 3, and a padding of 2. For an input feature map of shape $28 \times 28 \times 28$, what are the dimensions/shape of the output tensor after applying the Convolution Layer to the input?

0
1
2

2.7 You are given a Convolutional Layer with: number of input channels 3, number of filters 5, kernel size 4, stride 2, padding 1. What is the total number of trainable parameters for this layer? Don't forget to consider the bias.

0
1
2

2.8 You are given a fully-connected network with 2 hidden layers, the first of has 10 neurons, and the second hidden layer contains 5 neurons. Both layers use dropout with probability 0.5. The network classifies gray-scale images of size 8×8 pixels as one of 3 different classes. All neurons include a bias. Calculate the total number of trainable parameters in this network.

0
1
2

DO NOT SCAN/UPLOAD

2.5 Wrong name, but weights are forced to be zero
L₂ regularization with (6.1)

$$w^+ = w - \eta \triangleright (L(w) + \frac{1}{2} \lambda \|w\|_2^2)$$

$$w^+ = w - \eta \triangleright L - \eta \lambda w$$

$$\underline{w^+ = (1 - \eta \lambda) w - \eta \triangleright L}$$

$$\eta \lambda \ll 1 \quad \therefore \text{Each iteration pushes } w \text{ to 0}$$

2.6. $\frac{F + 2P - N}{S} + 1$

$$= \frac{28 + 2 \times 2 - 5}{3} + 1$$

$$= 10$$

$$10 \times 10 \times 6$$

2.7 (3, 5, 4, 4)

$$(3 \times 4 \times 4 + 1) \times 5$$

$$= 245$$

64.

2.8 8×8

10 5 8

0.5 0.5

$$64 \times 10 + 10 + 10 \times 5 + 5 + 5 \times 3 + 3$$

$$= 650 + 50 + 18$$

$$= \cancel{725} \quad 723$$

0  2.9 "*Breaking the symmetry*": Why is initializing all weights of a fully-connected layer to the same value problematic?

1

2

0  2.10 Explain the difference between *Auto-Encoders* and *Variational Auto-Encoders*.

1

2

0  2.11 Generative Adversarial Networks (GANs): What is the input to the generator network (1 pt)? What are the two inputs to the discriminator (1 pt)?

1

2

0  2.12 Explain how LSTM networks often outperform traditional RNNs. What in their architecture enables this?

1

2

0  2.13 Explain how batch normalization is applied differently between a fully connected layer and a convolutional layer (1 pt). How many learnable parameters does batch normalization contain following (a) a single fully-connected layer (1 pt), and (b) a single convolutional layer with 16 filters (1 pt)?

1

2

3

DO NOT SCAN/UPLOAD

2.9 Because all neuron have same weight
→ cause same gradients when Backpropagation
→ learn same things
→ NV update very slow or not optimal

2.10 Autoencoder: encoder encode the input feature to a low dimension space. Decoder decode from this space and use learned feature to generate output
VAE encode encode the input feature to a latent space which is Gaussian distributed, need sample the vector from the gaussian latent space to generate output

2.11 Real world Image / Random variable from latent space
and sample and sample from generator

2.12 Cell - way for Information transform and
flipping for gradient.

Also use tanh and sigmoid as Activation layer.
Also there are other gates to update information

2.13 Mini-Batch is normalized for all neuron in FC
while Mini-Batch is normalized for all
2 β and γ channels in conv layer

2 for ~~one~~ FC and 2x16 for conv

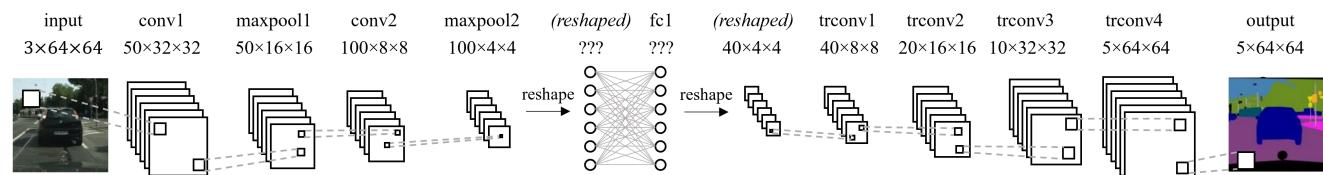
Problem 3 Convolutions (13 credits)

You are asked to perform **per-pixel** semantic segmentation on the Cityscapes dataset, which consists of RGB images of European city streets, and you want to segment the images into 5 classes (vehicle, road, sky, nature, other). You have designed the following network, as seen in the illustration below:

For clarification of notation: The shape **after** having applied the operation ‘conv1’ (the first convolutional layer in the network) is $50 \times 32 \times 32$.

You are using 2D convolutions with: `stride = 2` , `padding = 1` , and `kernel_size = 4` .

For the MaxPool operation, you are using: `stride = 2` , `padding = 0` , and `kernel_size = 2` .



3.1 What is the shape of the weight matrix of the fully-connected layer ‘fc1’? (Ignore the bias)

0
1
2

3.2 Explain the term ‘receptive field’ (1p). What is the receptive field of one pixel of the activation map. after performing the operation ‘maxpool1’(1p)? What is the receptive field of a single neuron in the output of layer ‘fc1’ (1p)?

0
1
2
3

DO NOT SCAN/UPLOAD

$$100 \times 4 \times 4 = 1600 \text{ Input}$$

$$40 \times 4 \times 4 = 640 \text{ Output}$$

Input shape 1600×640

How many pixels in the input image can affect only one pixel on the output image

$$\frac{F-N}{S} + 1 = 1, \quad F = 2$$

$$\frac{F-2}{2} = 0 \quad 2 \times 2$$

$$\frac{F+2P-N}{S} + 1 = 2$$

$$\frac{F+2-4}{2} = 2 \quad F = 4 \times 4$$

$$\frac{F-N}{S} + 1 = 2$$

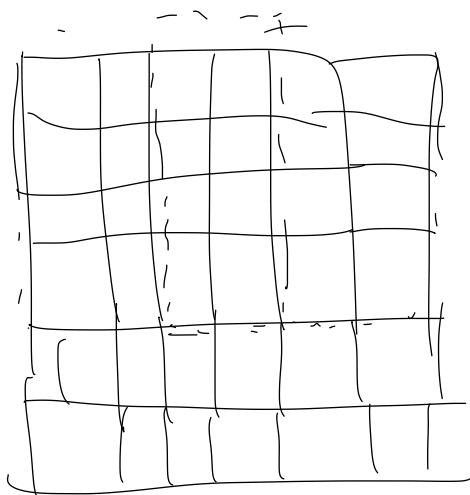
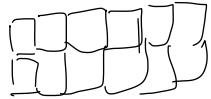
$$I = 6 \times 6$$

Whole Images, how all connected

$$(64 \times 64)$$

$$\frac{4+2-4}{2} + 1 = 2$$

$$(2 + 3 \times 2) \times 2$$



- 0 3.3 You now want to be able to classify finer-grained labels, which comprise of 30 classes. What is the **minimal** change in network architecture needed in order to support this without adding any additional layers?

- 0 3.4 Luckily, you found a pre-trained version of this network, which is trained on the original 5 labels. (It outputs a tensor of shape $5 \times 64 \times 64$). How can you make use of/build upon this pre-trained network (as a black-box) to perform segmentation into 30 classes.

- 0 3.5 Luckily, you have gained access to a large dataset of city street images. Unfortunately, these images are not labelled, and you do not have the resources to annotate them. However, how can you still make use of these images to improve your network? Explain the architecture of any networks that you will use and explain how training will be performed. (Note: This question is independent of (3.3) and (3.4))

- 0 3.6 Instead of taking 64×64 images as input, you now want to be able to train the network to segment images of arbitrary size > 64 . List, explicitly, two different approaches that would allow this. Your new network should support varying image sizes in run-time, without having to be re-trained.

3.3 Changing the last layer into FC

3.4 Use pretrained network and frozen the weight
only train the last FC layer

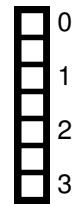
3.5 Use Autoencoder to learn the feature from
unlabeled images

And use pretrained encoder and classifier

3.6 Add a conv ~~at~~ first layer to deal the size
problem

Problem 4 Optimization (13 credits)

4.1 Explain the idea behind the RMSProp optimizer. How does it enable faster convergence than standard SGD? How does it make use of the gradient?

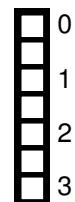


4.2 What is the *bias correction* in the ADAM optimizer? Explain the problem that it fixes.



softwares

4.3 You read that when training deeper networks, you may suffer from the *vanishing gradients* problem. Explain what are vanishing gradients in the context of deep convolutional networks and the underlying cause of the problem.



DO NOT SCAN/UPLOAD

① Use second Momentum to predict the next optimal factors and jump to it then make some correction, speed up the optimal process

RMSPROP is an adaptive learning rate method. It scales the learning rate enable fast converge by skipping the saddle point. Dampening the oscillation in the dim with high variance.

Second Momentum is a exponentially weighted moving average of square gradients

When accumulating gradient in weighted moving average fashion, the first gradient is initialized to 0, this biases all the accumulating gradient toward zero.

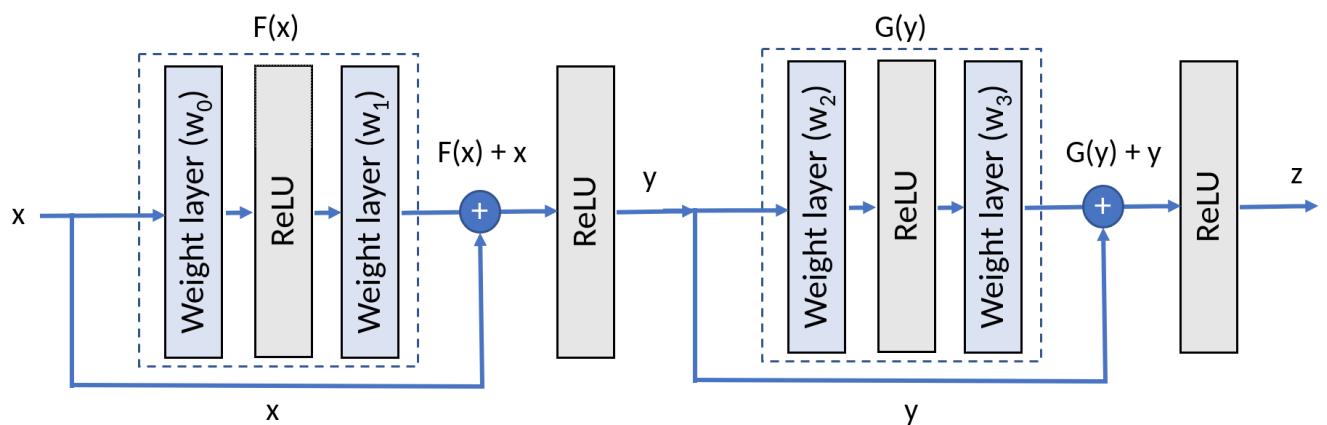
Solution: Bias correction normalized the accumulating the gradient in early step. $\sqrt{\text{magnitude of}}$

When the network is very deep, i.e. with a lot of layers, the gradient of deep layer are almost equal to zero, this situation will speed down the optimal process, which called gradient vanishing

- Activation layer saturates eg tanh sym
- Chain rule, the weights from early layer decay on last layer which also equal to 0

0
1
2
3
4
5

4.4 In the following image you can see a segment of a very deep architecture that uses residual connections. How are residual connections helpful against vanishing gradients? Demonstrate this mathematically by performing a weight update for w_0 . Make sure to explain how this reduces the effect of vanishing gradients. Hint: Write the mathematical expression for $\frac{\partial z}{\partial w_0}$ w.r.t all other weights.



DO NOT SCAN/UPLOAD

The Residual connection provide a highway for gradient, which can efficiently skip the gradient vanishing in deep backwork

Residual connection.

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

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

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

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

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

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

$$\begin{aligned} \frac{\partial z}{\partial w_0} &= \frac{\partial \text{ReLU}(G(y) + y)}{\partial (G(y) + y)} \cdot \frac{\partial G(y) + y}{\partial y} \\ &\quad - \frac{\partial \text{ReLU}(f(x) + x)}{\partial f(x) + x} \cdot \frac{\partial f(x) + x}{\partial w_0} \\ &= 1 \cdot \left(\frac{\partial w_3 \cdot \text{ReLU}(w_2 \cdot y)}{\partial \text{ReLU}(w_2 \cdot y)} \cdot \frac{\partial \text{ReLU}(w_2 \cdot y)}{\partial w_2 \cdot y} \cdot \frac{\partial w_2 \cdot y}{\partial y} + 1 \right) \\ &\quad - 1 \cdot \frac{\partial w_1 \cdot \text{ReLU}(w_0 \cdot x)}{\partial \text{ReLU}(w_0 \cdot x)} \cdot \frac{\partial \text{ReLU}(w_0 \cdot x)}{\partial w_0 \cdot x} \cdot \frac{\partial w_0 \cdot x}{\partial w_0} \end{aligned}$$

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

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

Compare to the gradient without skipconn:

there exists addition " $w_1 \cdot x$ " to parent gradient when

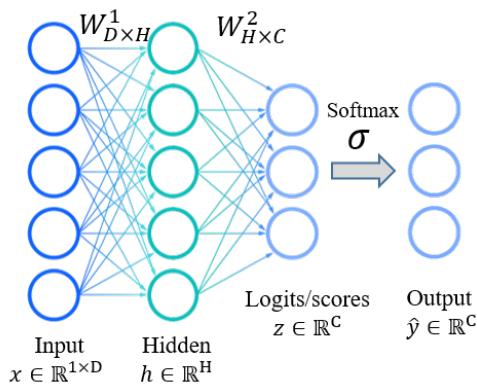
$$\text{e.g. } w_3 \rightarrow 0$$

Problem 5 Multi-Class Classification (18 credits)

Note: If you cannot solve a sub-question and need its answer for a calculation in following sub-questions, mark it as such and use a symbolic placeholder (i.e., the mathematical expression you could not explicitly calculate + a note that it is missing from the previous question.)

Assume you are given a labeled dataset $\{X, y\}$, where each sample x_i belongs to one of $C = 10$ classes. We denote its corresponding label $y_i \in \{1, \dots, 10\}$. In addition, you can assume each data sample is a row vector.

You are asked to train a classifier for this classification task, namely, a 2-layer fully-connected network. For a visualization of the setting, refer to the following illustration:



5.1 Why does one use a **Softmax** activation at the end of such a classification network? What property does it have that makes it a common choice for a classification task?

- 0
- 1
- 2

5.2 For a vector of logits \vec{z} , the Softmax function $\sigma : \mathbb{R}^C \rightarrow \mathbb{R}^C$, is defined:

$$\hat{y}_i = \sigma(\vec{z})_i = \frac{e^{z_i}}{\sum_{j=1}^C e^{z_j}}$$

where C is the number of classes and z_i is the i -th logit.

A special property of this function is that its derivative can be expressed in terms of the Softmax function itself. How could this be advantageous for training neural networks?

- 0
- 1

DO NOT SCAN/UPLOAD

softmax activation is to calculate the probability of every one from
the goal

last layer

softmax can compute the probability distribution of the output over the different classes

softmax function ensure all the output is valid probability and

sum up is equal to 1

This is easy to interpret the output as the prediction classes

In forward pass, you must calculate the probability distribution, in the same time you can cache it so that you can use it as the denominator

When Back propagation,

It sums up value from loss

0  5.3 Show explicitly how this can be done, by writing $\frac{\partial \hat{y}_i}{\partial z_i}$ in terms of \hat{y}_i .

1
2
3

0  5.4 Similarly, show explicitly how this can be done, by writing $\frac{\partial \hat{y}_i}{\partial z_j}$ in terms of \hat{y}_i and \hat{y}_j , for $i \neq j$.

1
2

DO NOT SCAN/UPLOAD

$i = j$

$$\begin{aligned}
 \frac{\partial \hat{y}_i}{\partial z_i} &= \frac{\partial g(\vec{z})_i}{\partial z_i} = \frac{e^{z_i} \cdot \sum_{j=1}^c e^{z_j} - e^{z_i} \cdot e^{z_i}}{\left(\sum_{j=1}^c e^{z_j} \right)^2} \\
 &= \frac{e^{z_i} \left(\sum_{j=1}^c e^{z_j} - e^{z_i} \right)}{\sum_{j=1}^c e^{z_j} \cdot \sum_{j=1}^c e^{z_j}} \\
 &= \frac{e^{z_i}}{\sum_{j=1}^c e^{z_j}} \cdot \left(1 - \frac{e^{z_i}}{\sum_{j=1}^c e^{z_j}} \right) \\
 &= \hat{y}_i (1 - \hat{y}_i)
 \end{aligned}$$

$i \neq j$

$$\begin{aligned}
 \frac{\partial \hat{y}_i}{\partial z_j} &= \frac{\partial g(\vec{z})_i}{\partial z_j} = \frac{0 \cdot \sum_{j=1}^c e^{z_j} - e^{z_i} \cdot e^{z_j}}{\left(\sum_{j=1}^c e^{z_j} \right) \left(\sum_{j=1}^c e^{z_j} \right)} \\
 &= - (\hat{y}_i) (\hat{y}_j) \\
 &= - \hat{y}_i \hat{y}_j
 \end{aligned}$$

5.5 Using the Softmax activation, what loss function $\mathcal{L}(y, \hat{y})$ would you want to *minimize*, to train a network on such a multi-class classification task? Name this loss function (1 pt), and write down its formula (2 pt), for a single sample x , in terms of the network's prediction \hat{y} and its true label y . Here, you can assume the label $y \in \{0, 1\}^C$ is a one-hot encoded vector:

$$y_i = \begin{cases} 1, & \text{if } i == \text{true class index} \\ 0, & \text{otherwise} \end{cases}$$

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

5.6 Having done a forward pass with our sample x , we will back-propagate through the network. We want to perform a gradient update for the weight $w_{j,k}^2$ (the weight which is in row j , column k of the second weights' matrix W^2). First, use the chain rule to write down the derivative $\frac{\partial \mathcal{L}}{\partial w_{j,k}}$ as a product of 3 partial derivatives (no need to compute them). For convenience, you can ignore the bias and omit the 2 superscript.

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

DO NOT SCAN/UPLOAD

Cross Entropy

$$L = - \sum_{j=1}^C y_j \cdot \ln \hat{y}_j$$

$$\frac{\partial L}{\partial w_{j,k}} = \frac{\partial L}{\partial \hat{y}_j} \cdot \frac{\partial \hat{y}_j}{\partial z_j} \cdot \frac{\partial z_j}{\partial w_{j,k}}$$

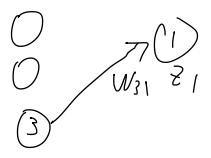
$$\frac{\partial L}{\partial w_{j,k}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial w_{j,k}}$$

0 5.7 Now, compute the gradient for the weight: $w_{3,1}^2$. For this, you will need to compute each
1 of the partial derivatives you have written above, and perform the multiplication to get the final
2 answer. You can assume the ground-truth label for the sample was `true_class = 3`. **Hint:** The
3 derivative of the logarithm is $(\log t)' = \frac{1}{t}$.

4

5

DO NOT SCAN/UPLOAD



$$\frac{\partial L}{\partial w_{j,k}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial w_{j,k}}$$

$$= -y_3 \cdot \frac{1}{\hat{y}_3} \cdot -\hat{y}_3 \cdot \hat{y}_1 \cdot h_3$$

$$= -1 \cdot 1 \cdot -1 \cdot \hat{y}_1 \cdot h_3$$

$$= \hat{y}_1 \cdot h_3$$

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

DO NOT SCAN/UPLOAD

DO NOT SCAN/UPLOAD