

Problem set 1

Deep Learning E1394

Out on Sept 18
Due on Oct 2, 23:59

Submit your written answers as a pdf typed in \LaTeX together with your code. Submit one answer per group (as assigned on Moodle) and include names of all group members in the document. Round answers to two decimal places as needed. Include references to any external sources you have consulted (points are deducted if those were used but not cited). See “Submission” at the bottom of the problem set for more details on how to submit using Github classroom.

1 Theoretical part

These are some recap and refresher problems to get you up to speed with the mathematics and statistical learning that is behind deep learning, as well as problems that help you work out how neural networks are trained.

1.1 Optimization

Let $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ be defined by $f(x, y) = x^2y + xy^2 - 6xy$. Identify and classify all critical points of f .

In order to identify the critical points, we first need to find the partial derivatives and set them equal to zero:

$$\begin{aligned} f : \mathbb{R}^2 &\rightarrow \mathbb{R}, f(x, y) = x^2y + xy^2 - 6xy \\ \frac{\partial f(x, y)}{\partial x} &= 2xy + y^2 - 6y = 0 \quad I \\ \frac{\partial f(x, y)}{\partial y} &= x^2 + 2xy - 6x = 0 \quad II \end{aligned}$$

Then, we can find all critical points of $f(x, y)$:

$$\text{I: } y(y + 2x - 6) = 0$$

$$\text{II: } x(x + 2y - 6) = 0$$

setting y to zero, we get in equation II

$$x(x - 6) = 0$$

setting this to zero, we arrive at $x_1 = 0$ and $x_2 = 6$

thus we have :

$$C_1 : (0, 0)$$

$$C_2 : (6, 0)$$

since the two equations mirror each other, we also have :

$$C_3 : (0, 6)$$

To find the last critical point, we solve the second term of the first equation (I_a) for y and plug it in the second part of the second equation (II_a). Solving for x we get the last critical point:

$$I_a : y + 2x - 6 = 0$$

$$y = 6 - 2x$$

$$II_a : x + 2(6 - 2x) - 6 = 0$$

$$x + 12 - 4x - 6 = 0$$

$$6 - 3x = 0 \Rightarrow 3x = 6 \quad x = 2$$

$$y = 6 - 2 \cdot 2 = 2$$

$$C_4 : (2, 2)$$

We put together the Hessian matrix, and use the determinant to classify the critical points:

$$H = \begin{bmatrix} \frac{\partial^2}{\partial x^2} & \frac{\partial}{\partial y \partial x} \\ \frac{\partial}{\partial x \partial y} & \frac{\partial^2}{\partial y^2} \end{bmatrix}$$

$$\begin{aligned}
\frac{\partial f(x,y)}{\partial x} &= 2xy + y^2 - 6y = 0 \\
\frac{\partial f(x,y)}{\partial y} &= x^2 + 2xy - 6x = 0 \\
\frac{\partial^2}{\partial x^2} &= 2y \\
\frac{\partial^2}{\partial y^2} &= 2x \\
\frac{\partial^2}{\partial x \partial y} &= \frac{\partial^2}{\partial y \partial x} = 2x + 2y - 6 \\
H_{f(x,y)} &= \begin{bmatrix} 2y & 2x + 2y - 6 \\ 2x + 2y - 6 & 2x \end{bmatrix} \\
\det(H_{f(x,y)}) &= 4xy - (2x + 2y - 6)^2 = D(x,y) \\
D(0,0) &= 0 - (-6)^2 = -36 < 0 \text{ saddle point} \\
D(6,0) &= 0 - (2 \cdot 6 - 6)^2 \\
&= 0 - (6)^2 \\
&= -36 < 0 \text{ saddle point} \\
D(0,6) &= -36 < 0 = \text{ saddle point} \\
D(2,2) &= 4 \cdot 2 \cdot 2 - (2 \cdot 2 + 2 \cdot 2 - 6)^2 \\
&= 16 - 2^2 = 12 > 0 \text{ local minimum}
\end{aligned}$$

We found four critical points: $D(0,0)$, $D(6,0)$ and $D(0,6)$ are saddle points of the function, $D(2,2)$ is a local minimum of the function.

1.2 Activation functions

Compute the gradient of the function

$$f(b, w) = \text{ReLU}(b + xw) = \begin{cases} 0 & \text{if } b + xw < 0 \\ b + xw & \text{if } b + xw \geq 0 \end{cases}$$

with respect to the parameters b and w , with $x, b, w \in \mathbb{R}$.

generally we apply the chain rule:

with (i) $\partial \text{ReLU} = \partial \max(0, f(z))$

where $z = b + xw$

$$(ii) \frac{\partial f}{\partial b} \text{ and } \frac{\partial f}{\partial w} \Rightarrow \nabla f(b, w) = \begin{pmatrix} \frac{\partial f}{\partial b} \\ \frac{\partial f}{\partial w} \end{pmatrix}$$

this makes for 3 cases:

1. $b + xw < 0$, $\partial \text{ReLU}(f) = 0$ thus $0 \cdot \nabla f(b, w)0 = 0$
2. $b + xw > 0$

$$\begin{aligned}\frac{\partial \text{ReLU}(f)}{\partial f} &= 1 \\ \frac{\partial f(b, w)}{\partial b} &= 1 \\ \frac{\partial f(b, w)}{\partial w} &= x \\ \text{thus, } \nabla \left(f(b, w) \right) &= \begin{pmatrix} 1 \\ x \end{pmatrix} \\ \Rightarrow 1 \cdot \begin{pmatrix} 1 \\ x \end{pmatrix} &= \begin{pmatrix} 1 \\ x \end{pmatrix}\end{aligned}$$

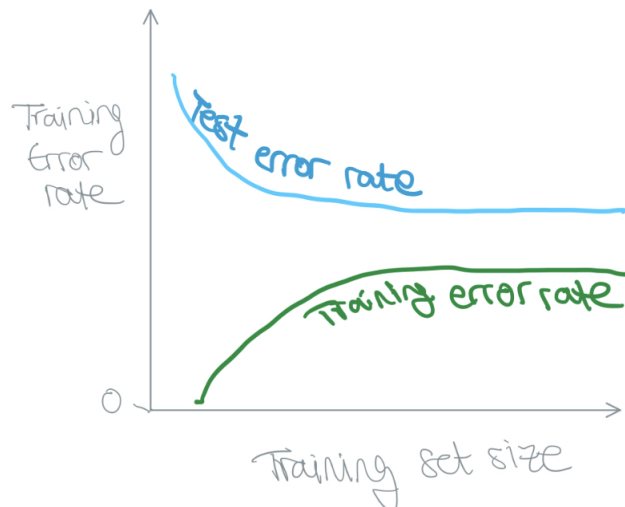
3. if $b + wx = 0$, the derivative of $\text{ReLU}(f)$ is undefined, since the left-hand limit (0) is different from the right-hand limit (1)

1.3 Overfitting

This question will test your general understanding of overfitting as it relates to model complexity and training set size. Consider a continuous domain and a smooth joint distribution over inputs and outputs, so that no test or training case is ever duplicated exactly. *The graphs in your solutions should be drawn on paper and then included as a picture. In your answers describe how your graphs should be read.*

1.3.1 Error rate versus dataset size

Sketch a graph of the typical behavior of training error rate (y-axis) as a function of training set size (x-axis). Add to this graph a curve showing the typical behavior of test error rate versus training set size, on the same axes. (Assume that we have an infinite test set drawn independently from the same joint distribution as the training set). Indicate on your y-axis where zero error is and draw your graphs with increasing error upwards and increasing training set size rightwards.



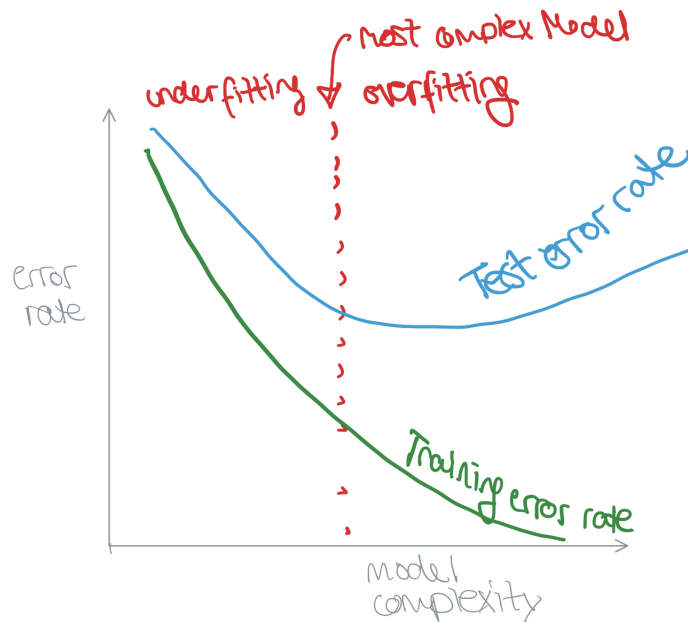
For very small training data sets, e.g., less than 10 observations, a typical model should achieve close to 0 training error, as it commands sufficient degrees of freedom to perfectly "map" onto the training data. This overfitting on the training data in turn is associated with a high testing error, since the model poorly generalizes to new data.

As the training set size increases, this counters the overfitting on the training data, making the training error go up. At the same time, this implies a decrease in the test error, as the model starts to better generalize.

Beyond a certain training test size, both training and test error rates usually stabilize.

1.3.2 Error rate versus model complexity

For a fixed training set size, sketch a graph of the typical behavior of training error rate (y-axis) versus model complexity (x-axis). Add to this graph a curve showing the typical behavior of the corresponding test error rate versus model complexity, on the same axes (again on an IID infinite test set). Show where on the x-axis you think is the most complex model that your data supports (mark this with a vertical line). Choose the x-range so that this line is neither on the extreme left nor on the extreme right. Indicate on your vertical axis where zero error is and draw your graphs with increasing error upwards and increasing complexity rightwards.



This graph illustrates the bias-variance trade-off: As model complexity increases, training error decreases (lower bias), but test error likely increases (higher variance) due to overfitting, resulting in the U-shaped test error curve.

The vertical dashed line highlights the optimal degree of complexity which is associated with the lowest test error rate.

1.3.3 Training epochs

Use a similar graph to illustrate the error rate as a function of training epochs in neural networks. One of the commonly used regularization methods in neural networks is early stopping. Describe (also using the graph) how early stopping is applied, and argue qualitatively why (or why not) early stopping is a reasonable regularization metric.



Typically, the validation loss is used for early stopping in Machine Learning, however in an idealized case, the validation loss should behave very similar to the test error.

As displayed in the graph, early stopping aims to detect the "plateau" in the test error, before overfitting begins and the test error would start to increase again.

In practice, the challenge of early stopping is to define the right "patience", meaning the number of epochs without improvement in the loss/error. If patience is too high, we would move beyond the optimum test error, if patience is too low, we might get stuck in a local minimum.

1.4 Capacity of neural network

We have a single hidden-layer neural network with two hidden units, and data with 2 features $X = (x_1, x_2)$. Design an activation function $g(\cdot)$ applied in the hidden layer, and an output activation function $o(\cdot)$ applied in the output layer, as well as a set of values of the weights that will create a logistic regression model where the input to the logistic function is of the form: $\beta_0 + \beta_1 x_1 x_2$.

$$f(x) = o \left(b^{[2]} + \sum_i^2 w_i^{[2]} \cdot g \left(b_i^{[1]} + \sum_j^2 w_{ij}^{[1]} x_j \right) \right)$$

$$g(\cdot) = (a(x_1, x_2))^2 = \text{quadratic input activation}$$

Proposing weight and bias terms for input activation:

$$b_1^{[1]} = 0, w_{11} = 1, w_{12} = 1$$

$$b_2^{[1]} = 0, w_{21} = 1, w_{22} = -1$$

$$h_1(x_1, x_2) = 0 + x_1 + x_2$$

$$h_2(x_1, x_2) = 0 + x_1 - x_2$$

$o(\cdot) = \sigma(\cdot)$ = sigmoid output activation to emulate a logistic regression

Proposing weight and bias terms for output activation:

$$b^{[2]} = \beta_0, w_1^{[2]} = \frac{1}{4}\beta_1, w_2^{[2]} = -\frac{1}{4}\beta_1$$

putting it together :

$$\begin{aligned} &= \frac{1}{4}\beta_1 (x_1 + x_2)^2 - \frac{1}{4}\beta_1 (x_1 - x_2)^2 \\ &= \frac{1}{4}\beta_1 [x_1^2 + 2x_1x_2 + x_2^2 - x_1^2 + 2x_1x_2 - x_2^2] \\ &= \frac{1}{4}\beta_1 [4x_1x_2] = \beta_1 x_1x_2 \\ &= \sigma(\beta_0 + \beta_1 x_1x_2) \\ &= \frac{1}{1 + e^{-z}}, \text{ with } z = \beta_0 + \beta_1 x_1x_2 \end{aligned}$$

1.5 Neural network theory

Derive the weight updates in gradient descent for a neural network with 2 hidden layers (superscripts [1] and [2]) that each have $H^{[1]}$ and $H^{[2]}$ hidden units respectively. The output layer has superscript [3], and we want to classify the data into K classes.

The output of the network for classes $k = 1, 2, \dots, K$ and one data point

$x \in \mathbb{R}^d$ can be written as such:

$$f_k(\theta; X) = o(a_k^{[3]}) = o(b_k^{[3]} + \sum_{m=1}^{H^{[2]}} w_{km}^{[3]} h_m^{[2]}) \quad (1)$$

$$h_m^{[2]} = \sigma(a_m^{[2]}) = \sigma\left(b_m^{[2]} + \sum_{i=1}^{H^{[1]}} w_{mi}^{[2]} h_i^{[1]}\right) \quad (2)$$

$$h_i^{[1]} = \sigma(a_i^{[1]}) = \sigma\left(b_i^{[1]} + \sum_{j=1}^d w_{ij}^{[1]} x_j\right), \quad (3)$$

where θ indicates the vector of all weights and biases.

While this is typically not recommended, in this exercise we use a quadratic loss function for classification. We use a softmax activation function at the output layer and a ReLU activation function at the hidden layers. Derive the weight update for the weights of the first hidden layer $w_{ij}^{[1]}$. Start by stating the gradient descent weight update for the $(r+1)^{th}$ iteration as a function of the $(r)^{th}$ iteration, and then compute the partial derivative needed in the update.

Show all the steps in your derivation. You may use the derivatives for activation functions introduced in class. There is no need to show the derivations of those.

$$w_{ij}^{[1](r+1)} = w_{ij}^{[1](r)} - \eta \nabla \left(w_{ij}^{[1](r)} \right)$$

$$\frac{\partial L(f_k(\theta, X))}{\partial w_{ij}^{[1]}} = \frac{\partial L}{\partial f_k} \cdot \frac{\partial f_k}{\partial a_k^{[3]}} \cdot \frac{\partial a_k^{[3]}}{\partial h_m^{[2]}} \cdot \frac{\partial h_m^{[2]}}{\partial a_m^{[2]}} \cdot \frac{\partial a_m^{[2]}}{\partial h_i^{[1]}} \cdot \frac{\partial h_i^{[1]}}{\partial a_i^{[1]}} \cdot \frac{\partial a_i^{[1]}}{\partial w_{ij}^{[1]}}$$

we then compute the derivatives for all these terms separately

Term 1 (derivative of loss)

$$\frac{\partial L}{\partial f_k} = \frac{\partial}{\partial f_k} (y - f_k(\theta, x))^2 = -2(y - f_k)$$

Term 2 (derivative of softmax activation function):

$$\begin{aligned} \frac{\partial f_k}{\partial a^{[3]}_k} \\ o(a^{[3]}) &= \frac{e^{a^{[3]}_k}}{\sum_{m=1}^k e^{a^{[3]}_m}} a^{[3]}_m \\ \frac{\partial o(a^{[3]}_k)}{\partial a^{[3]}_k} &= o(a^{[3]}_k) \left(\delta_{kk} - o(a^{[3]}_k) \right) \end{aligned}$$

using the Kronecker Delta

Term 3 (derivative of output activation)

$$\frac{\partial a^{[3]}_k}{\partial h^{[2]}_m} = \frac{\partial}{\partial h^{[2]}_m} \cdot \left[b^{[3]}_k + \sum_{m=1}^{H^{[2]}} w^{[3]}_{km} h^{[2]}_m \right] = w^{[3]}_{km}$$

Term 4 (second hidden layer activation function)

$$\frac{\partial h^{[2]}_m}{\partial a^{[2]}_m} = \sigma(a^{[2]}_m) = \text{ReLU}(a^{[2]}_m)$$

$$\frac{\partial}{\partial a^{[2]}_m} \text{ReLU}(a^{[2]}_m) = 1 \text{ if } a^{[2]}_m > 0, \text{ undefined if } a^{[2]}_m = 0 \text{ and } 0 \text{ if } a^{[2]}_m < 0$$

Term 5 (derivative of second hidden layer activation):

$$\frac{\partial a^{[2]}_m}{\partial h^{[1]}_i} = \frac{\partial}{\partial h^{[1]}_i} \cdot \left[b^{[2]}_m + \sum_{i=1}^{H^{[1]}} w^{[2]}_{mi} h^{[1]}_i \right] = w^{[2]}_{mi}$$

Term 6 (derivative of the first hidden layer activation function)

$$\frac{\partial h^{[1]}_i}{\partial a^{[1]}_i} = \sigma(a^{[1]}_i) = \text{ReLU}(a^{[1]}_i) = 1 \text{ if } a^{[1]}_i > 0, \text{ undefined if } a^{[1]}_i = 0 \text{ and } 0 \text{ if } a^{[1]}_i < 0$$

Term 7 (derivative of the first hidden layer activation)

$$\frac{\partial a^{[1]}_i}{\partial w^{[1]}_{ij}} = \frac{\partial}{\partial w^{[1]}_{ij}} \cdot \left[b^{[1]}_i + \sum_{j=1}^d w^{[1]}_{ij} x_j \right]$$

$$= x_j$$

Putting all terms together:

$$\frac{\partial L(fk(\theta, X))}{\partial w^{[1]}_{ij}} = -2(y - f_k) o(a^{[3]}_k) \left(\delta_{km} - o(a^{[3]}_k) w^{[3]}_{km} w^{[2]}_{mi} x_j \right)$$

$$\text{assumption : } a^{[1]}_i > 0 \text{ and } a^{[2]}_m > 0$$

$$\text{if either } a^{[1]}_i \text{ or } a^{[2]}_m < 0 \text{ then } \frac{\partial L(fk(\theta, X))}{\partial w^{[1]}_{ij}} = 0$$

$$\text{if either } a^{[1]}_i \text{ or } a^{[2]}_m = 0 \text{ then } \frac{\partial L(fk(\theta, X))}{\partial w^{[1]}_{ij}} = \text{undefined}$$

2 Neural network implementation

You will implement different classes representing a fully connected neural network for image classification problems. There are two classes of neural networks: one using only basic packages and one using PyTorch. In addition to that, a third class of neural network has been implemented using Tensorflow and requires some fixing in order to function correctly.

As you work through this problem, you will see how those machine learning libraries abstract away implementation details allowing for fast and simple construction of deep neural networks.

For each approach, a Python template is supplied that you will need to complete with the missing methods. Feel free to play around with the rest, but please do not change anything else for the submission. All approaches will solve the same classification task, which will help you validate that your code is working and that the network is training properly.

For this problem, you will work with the MNIST dataset of handwritten digits, which has been widely used for training image classification models. You will build models for a multiclass classification task, where the goal is to predict what digit is written in an image (to be precise, this is a k -class classification task where in this case $k = 10$). The MNIST dataset consists of black and white images of digits from 0 to 9 with a pixel resolution of 28x28. Therefore, in a tensor representation the images have the shape 28x28x1. The goal is to classify what digit is drawn on a picture using a neural network with the following characteristics:

- an arbitrary amount of hidden layers, each with arbitrary amount of neurons
- sigmoid activation function for all hidden layers
- softmax activation function for the output layer
- cross entropy loss function. *

** For the implementation from scratch (Part (a)) we use a mean squared error (MSE) loss function. This is not recommended for a classification task, but we use it to simplify the implementation. In the future please consider using cross entropy / log loss instead.*


../img/MNIST-few-examples.png

Figure 1: Example images from MNIST dataset with 28x28 pixel resolution.

The training set with n data points is of the form

$$(x^{(i)}, y^{(i)}), \quad i = 1, \dots, n, \quad (4)$$

with $y^{(i)} \in \{0, 1\}^k$, where $y_j^{(i)} = 1$ when j is the target class, and 0 otherwise (i.e., the output is one-hot encoded). The softmax output function $\hat{y} = h_\theta(x)$ outputs vectors in \mathbb{R}^k , where the relative size of the \hat{y}_j corresponds roughly to how likely we believe that the output is really in class j . If we want the model to predict a single class label for the output, we simply predict class j for which \hat{y}_j takes on the largest value.

Tasks

1. Complete the implementation of the neural network classes marked by TODO comments:
 - (a) From scratch (`network_scratch.py`)
 - (b) Using PyTorch (`network_pytorch.py`)

A repository with the template files will automatically be created for your team when registering for the GitHub classroom assignment at https://classroom.github.com/a/XgELWE_C.
2. In `network_tensorflow.py` a Neural Network class has been implemented using Tensorflow and Keras.
 - (a) There are three mistakes in the code preventing the network from working correctly. Comment the mistakes out and repair the class.
 - (b) Implement a time-based learning rate class `TimeBasedLearningRate`: it should be initialized with a positive integer as initial learning rate and have the learning rate reduced by 1 at each step, until a learning rate of 1 is reached.
3. In addition to the implementation, show how your networks perform on the MNIST classification task in `MNIST_classification.ipynb`.
 - (a) Plot the accuracy of each neural network on the training and on the validation set as a function of the epochs
 - (b) Please make sure that the uploaded Notebook shows printed training and validation progress.
 - (c) Please add team number, team member names and matriculation numbers as a comment at the top of the notebook.

Please note that the classifier does not need to get anywhere close to 100% accuracy. Especially, with the small amount of training data this is difficult. Instead, the training output should indicate that the model learns something, i.e., the accuracy increases over the epochs on the training data and is significantly better than random guessing.

Submission

The submission of the whole problem set is done via GitHub classroom. You have to register for the assignment with your GitHub account at https://classroom.github.com/a/XgELWE_C. Please make sure that your GitHub account profile includes your real name. When starting the assignment, please create or join a team according to the teams assigned in Moodle (use the exact team name from Moodle).

Please upload / push all solutions to the GitHub repository which was created for your team. Please upload your solutions for the theoretical part (1) as a PDF to GitHub. If you upload multiple PDF files, please indicate in the file name to which subtask they are corresponding (we much prefer one single pdf). For the practical part (2), just push your code changes to the existing files.

Anybody in your team can push as often as they want. Once the deadline has passed, you are not able to push to your repository anymore and all changes on your **main** branch will be considered for grading. See the README.md in the repository for more details on how to push your changes to GitHub.