



ECE 046211 - Technion - Deep Learning

HW2 - Multilayer NNs and Convolutional NNs



Keyboard Shortcuts

- Run current cell: **Ctrl + Enter**
- Run current cell and move to the next: **Shift + Enter**
- Show lines in a code cell: **Esc + L**
- View function documentation: **Shift + Tab** inside the parenthesis or `help(name_of_module)`
- New cell below: **Esc + B**
- Delete cell: **Esc + D, D** (two D's)



Students Information

- Fill in

Name	Campus Email	ID
Student 1	student_1@campus.technion.ac.il	123456789
Student 2	student_2@campus.technion.ac.il	987654321



Submission Guidelines

- Maximal grade: 100.
- Submission only in **pairs**.
 - Please make sure you have registered your group in Moodle (there is a group creation component on the Moodle where you need to create your group and assign members).
- **No handwritten submissions.** You can choose whether to answer in a Markdown cell in this notebook or attach a PDF with your answers.

- **SAVE THE NOTEBOOKS WITH THE OUTPUT, CODE CELLS THAT WERE NOT RUN WILL NOT GET ANY POINTS!**
- What you have to submit:
 - If you have answered the questions in the notebook, you should submit this file only, with the name: `ece046211_hw2_id1_id2.ipynb`.
 - If you answered the questions in a different file you should submit a `.zip` file with the name `ece046211_hw2_id1_id2.zip` with content:
 - `ece046211_hw2_id1_id2.ipynb` - the code tasks.
 - `ece046211_hw2_id1_id2.pdf` - answers to questions.
 - No other file-types (`.py` , `.docx` ...) will be accepted.
- Submission on the course website (Moodle).
- **Latex in Colab** - in some cases, Latex equations may no be rendered. To avoid this, make sure to not use *bullets* in your answers ("* some text here with Latex equations" -> "some text here with Latex equations").



Working Online and Locally

- You can choose your working environment:
 1. Jupyter Notebook , **locally** with [Anaconda](#) or **online** on [Google Colab](#)
 - Colab also supports running code on GPU, so if you don't have one, Colab is the way to go. To enable GPU on Colab, in the menu: `Runtime` → `Change Runtime Type` → `GPU`.
 2. Python IDE such as [PyCharm](#) or [Visual Studio Code](#).
 - Both allow editing and running Jupyter Notebooks.
- Please refer to [Setting Up the Working Environment.pdf](#) on the Moodle or our GitHub (<https://github.com/taldatech/ee046211-deep-learning>) to help you get everything installed.
- If you need any technical assistance, please go to our Piazza forum (`hw2` folder) and describe your problem (preferably with images).



Agenda

- [Part 1 - Theory](#)
 - [Q1 - Generalization in A Teacher-Student Setup](#)
 - [Q2 - "Typical" Generalization in Multilayer Neural Networks](#)
 - [Q3 - Deep Double Descent](#)
 - [Q4 - Initialization](#)
 - [Q5 - Invariance and Equivariance](#)
 - [Q6 - VGG Architecture](#)
- [Part 2 - Code Assignments](#)

- Task 1 - The Importance of Activation and Initialization
- Task 2 - MLP-based Deep Classifier
- Task 3 - Design a CNN
- Credits



Part 1 - Theory

- You can choose whether to answer these straight in the notebook (Markdown + Latex) or use another editor (Word, LyX, Latex, Overleaf...) and submit an additional PDF file, **but no handwritten submissions**.
- You can attach additional figures (drawings, graphs,...) in a separate PDF file, just make sure to refer to them in your answers.
- *L^AT_EX* [Cheat-Sheet](#) (to write equations)
 - [Another Cheat-Sheet](#)



Question 1 - Generalization in A Teacher-Student Setup

Recall from lecture 4 the Bayes Risk $\overline{\mathcal{R}}(w)$:

$$\overline{\mathcal{R}}(w) \triangleq \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \sigma_\epsilon^2 I), w_{true} \sim \mathcal{N}(0, \frac{\sigma_w^2}{d} I)} [\mathcal{R}],$$

where,

$$\mathcal{R}(w_\mu) = \|w_\mu - w_{true}\|^2 = \|(H_\mu^{-1} H - I)w_{true} + H_\mu^{-1} X^T \epsilon\|^2$$

Prove:

$$\overline{\mathcal{R}}(w_\mu) = \sum_{i=1}^d \frac{(\sigma_w^2/d)\mu^2 + \sigma_\epsilon^2 \lambda_i}{(\lambda_i + \mu)^2}$$

Hints:

- $\mathbb{E} [\epsilon^T X H_\mu^{-1} H_\mu^{-1} X^T \epsilon] = \sum_{i,j}^N \mathbb{E}[\epsilon_i \epsilon_j] (X H_\mu^{-1})_i (H_\mu^{-1} X^T)_j$
- $\mathbb{E}[\epsilon_i \epsilon_j] = \sigma_\epsilon^2 \delta_{ij}$
- $\sum_{i=1}^N (X H_\mu^{-1})_i (H_\mu^{-1} X^T)_i = \text{Tr} [X H_\mu^{-2} X^T]$



Question 2 - "Typical" Generalization in Multilayer Neural Networks

We examine a "student" neural network $f_{\mathbf{w}}(\mathbf{x})$ with parameter vector $\mathbf{w} \in \mathbb{R}^k$ and input $\mathbf{x} \in \mathbb{R}^{d_0}$ used in a binary classification problem where the training set is $\mathcal{S} = \{\mathbf{x}^{(n)}\}_{n=1}^N$ sampled i.i.d. from P_X , where the binary (± 1) labels are generated by a "teacher" neural network $f_{\mathbf{w}_*}(\mathbf{x})$ with the same architecture. To understand the "typical" generalization of the student, we examine the following "Guess and Check" algorithm to learn its weights: we randomly sample parameters vectors $\mathbf{w}_1, \mathbf{w}_2, \dots$ i.i.d. from P_W , in which each parameter is sampled independently from a uniform distribution over $Q = \{-(q-1)/2, \dots, -1, 0, 1, \dots, (q-1)/2\}$ quantization levels, where $q = |Q|$ is an odd positive number (assume that teacher weights are also in Q). We do this until a stopping time t in which we perfectly fit the dataset: $\forall n : f_{\mathbf{w}_t}(\mathbf{x}^{(n)}) = f_{\mathbf{w}_*}(\mathbf{x}^{(n)})$. We examine a two-layer neural network with d_1 hidden neurons

$$f_{\mathbf{w}}(\mathbf{x}) = \text{sign}(\mathbf{w}_2^\top [\mathbf{W}_1 \mathbf{x}]_+)$$

where $[\cdot]_+$ is the ReLU activation function, the teacher has at most $d_1^* < d_1$ non-zero neurons (i.e., the other $d_1 - d_1^*$ hidden neurons in the teacher to have all the incoming and outgoing weights equal to zero). Each of the teacher's weights are also in Q .

1. Calculate the probability $P_{\mathbf{w} \sim P_W}(\mathbf{w} = \mathbf{w}_*)$.
2. Prove that

$$(1) \quad p_* \triangleq P_{\mathbf{w} \sim P_W}(\forall \mathbf{x} : f_{\mathbf{w}}(\mathbf{x}) = f_{\mathbf{w}_*}(\mathbf{x})) \geq q^{-d_0 d_1^* - d_1}. \quad (1)$$

3. Show that for any constant $T > 0$, we can bound the probability of stopping time $t > T$ as

$$(2) \quad [T] \leq \frac{\log P(t > T)}{\log(1 - p_*)}. \quad (2)$$

4. Prove the generalization bound:

Theorem 1 *With probability* $(1 - \eta)(1 - \delta)$,

$$(3) \quad \epsilon < \frac{(d_0 d_1^* + d_1) \log q + \log \frac{1}{\delta} + \log \log \frac{1}{\eta}}{N} \quad (3)$$

Hint: Combine the results from previous sections, using the approximations $[T] \approx T$ and $\log(1 - p_*) \approx -p_*$ (treat these approximations as exact), and the following basic generalization bound (which we learned in class):

Theorem 2 *For any* $f \in \mathcal{F}$ $f \in \mathcal{F}$, *with probability* $1 - \delta$,

$$(4) \epsilon \triangleq \mathbb{P}_{\mathbf{x}} (f_{\mathbf{w}}(\mathbf{x}) \neq f_{\mathbf{w}_*}(\mathbf{x})) < \frac{\log|\mathcal{F}| + \log \frac{1}{\delta}}{N}. \quad (4)$$

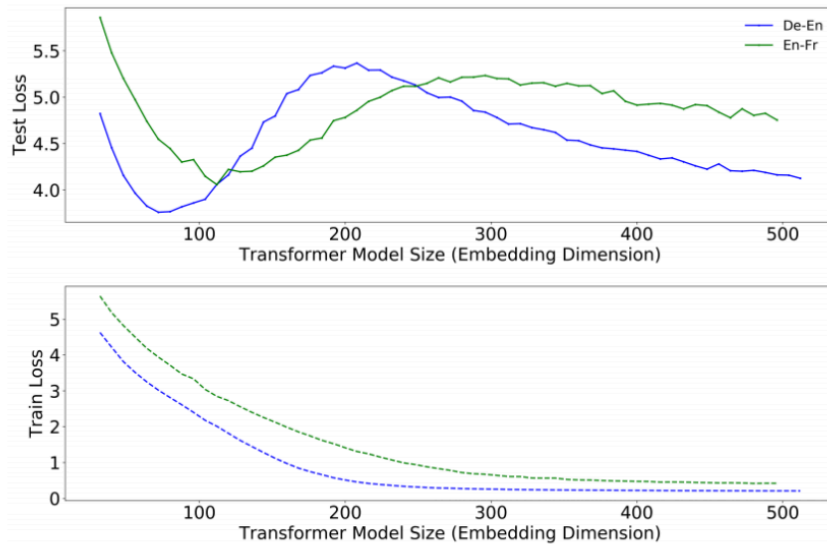
6. Is the bound in eq. (3) better than the bound in eq. (4) in which $\mathcal{F} = \{f_{\mathbf{w}} : \mathbf{w} \in Q^k\}$ is the student hypothesis class (in which each parameter can have one of q values)? Explain and ignore the (negligible) $\log \log \frac{1}{\eta}$ term.



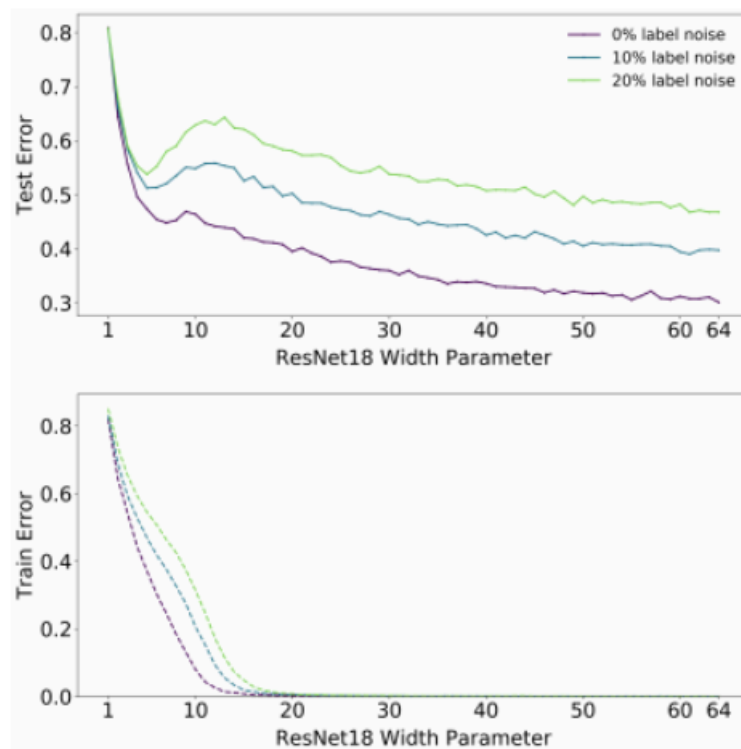
Question 3 - Deep Double Descent

For the following plots:

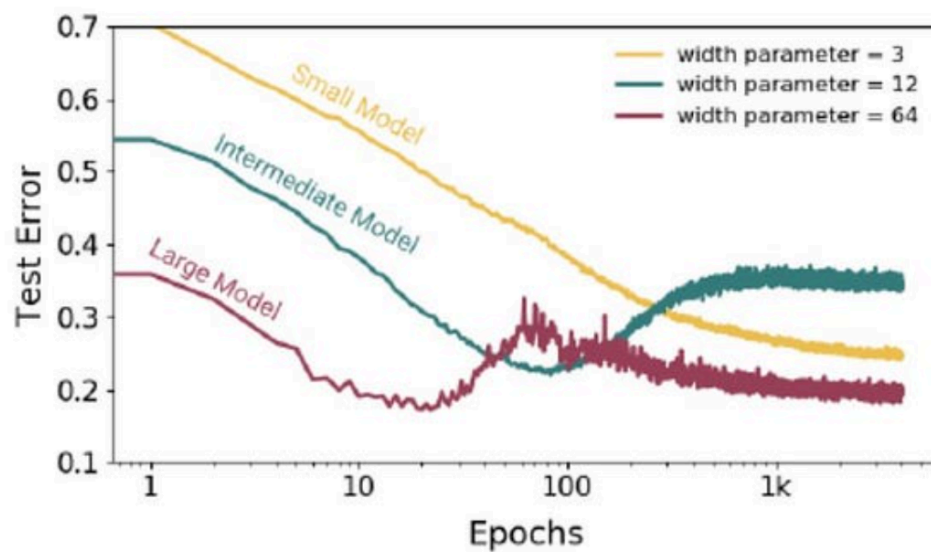
1. Where is the critical point (the point of transition between the "Classical Regime" and "Modern Regime") of the deep double descent?
2. What type of double descent is shown (**look closely at the graph**)? Explain. There can be more than one correct answer.



a.



b.



c.



Question 4 - Initialization

Recall that in lecture 5 we were discussing how to calculate the initialization variance, and reached the conclusion that

$$\sigma_l = \frac{1}{\sqrt{\sum_j \mathbb{E}[\varphi^2(u_{l-1}[j])]}}$$

Show that for ReLU activation ($\varphi(z) = \max(0, z)$), the optimal variance satisfies:

$$\sigma_l = \sqrt{\frac{2}{d_{l-1}}}$$

1. Under the assumption that the distribution of W is symmetric (\rightarrow the distribution of u is symmetric).
2. Using the central limit theorem for large width.

Answer each section **separately** and assume the sections are independent.

All the notations are the same as in the lecture slides.



Question 5 - Invariance and Equivariance

Consider a linear regression problem with a neural network of the form $\hat{y} = x^T w$, where $x, w \in \mathbb{R}^d$. Assume the network is trained to minimize

$$\mathcal{L}_0(w) = \frac{1}{2} \sum_{n=1}^N \left(y^{(n)} - w^T x^{(n)} \right)^2 = \frac{1}{2} \|y - Xw\|^2,$$

where $X = [x^{(1)}, \dots, x^{(N)}]^T$ is the sample set and $y = [y^{(1)}, \dots, y^{(N)}]^T$ is the corresponding label vector.

We define operator G as some operator on the input x , e.g. the input is an image and G spins it clockwise 90° . G is defined via an *invertible* square matrix of dimension d . We define a function f as **invariant** w.r.t. G if $f(Gx) = f(x)$, i.e. applying G does not change the output of f . We will mark \mathcal{G} a set of operators, and consider functions that are invariant w.r.t. all operators in \mathcal{G} . For example, if \mathcal{G} is a set of clockwise rotation operators, an invariant f might be the sum of all pixel values.

The following are a set of properties that apply to rotation matrices (operators):

1. If $G \in \mathcal{G}$, $G^T \in \mathcal{G}$.
2. Any $G \in \mathcal{G}$ is *orthogonal* ($GG^T = G^T G = I$).
3. Applying operator $G \in \mathcal{G}$ on all elements of \mathcal{G} returns \mathcal{G} , meaning $\bigcup_{G' \in \mathcal{G}} GG' = \mathcal{G}$.

We define the following cost function

$$\mathcal{L}_1(w) = \frac{1}{2} \sum_{G \in \mathcal{G}} \|y - XG^T w\|^2$$

- a. Explain intuitively why this cost function encourages the predictor w to be invariant.
- b. Show that $\mathcal{L}_1(w)$ is invariant to rotation transformations of the form $w \rightarrow Gw$, for all $G \in \mathcal{G}$.
- c. Prove that the optimal w_1^* that minimizes $\mathcal{L}_1(w)$ is

$$w_1^* = K^{-1} R X^T Y,$$

where $K = \sum_{G \in \mathcal{G}} G^T X^T X G$, $R = \sum_{G \in \mathcal{G}} G$ and we assume K is invertible.

d. Show that $GR = R$ for all $G \in \mathcal{G}$.

e. Show that $GK^{-1}G^T = K^{-1}$ for all $G \in \mathcal{G}$.

Hint: for invertible A, B , $(AB)^{-1} = B^{-1}A^{-1}$

f. Prove that the optimal w_1^* yields a predictor $\hat{y}(x) = x^T w_1^*$ that is invariant to $x \rightarrow Gx$ for all $G \in \mathcal{G}$.

Suppose that the network is now of the form $\hat{y} = Wx$, where $x \in \mathbb{R}^d$, $W \in \mathbb{R}^{d \times d}$, $\hat{y} \in \mathbb{R}^d$. A function f is **equivariant** to operator G if $f(Gx) = Gf(x)$, i.e. applying the operator on the input is the same as applying the operator on the output. Let $X = [x^{(1)}, \dots, x^{(N)}]^T$ be the sample set and $Y = [y^{(1)}, \dots, y^{(N)}]^T$ is the corresponding label matrix.

g. Suggest a cost function $\mathcal{L}_2(w)$ that encourages the predictor $\hat{y} = Wx$ is equivariant for all $G \in \mathcal{G}$. Provide a short explanation.



Question 6 - VGG Architecture

1. The VGG-11 CNN architecture consists of 11 convolution (CONV)/fully-connected (FC) layers (every CONV layer has the same padding and stride, every MAXPOOL layer is 2×2 and has padding of 0 and stride 2). Fill in the table. You need to **consider the bias**.

- CONV M - N : a convolutional layer of size $M \times M \times N$, where M is the kernel size and N is the number of filters. *stride* = 1, *padding* = 1.
- POOL2: 2×2 Max Pooling with *stride* = 2
 - In case the input of the layer is odd, you should round down. For example, if the output of the layer should be $3.5 \times 3.5 \times 3$, you should round to $3 \times 3 \times 3$ (i.e., ignore the last column of the input image when performing MaxPooling).
- FC- N : a fully connected layer with N neurons.

Layer	Output Dimension	Number of Parameters (Weights)
INPUT	224x224x3	0
CONV3-64	-	-
ReLU	-	-
POOL2	-	-
CONV3-128	-	-
ReLU	-	-
POOL2	-	-

Layer	Output Dimension	Number of Parameters (Weights)
CONV3-256	-	-
ReLU	-	-
CONV3-256	-	-
ReLU	-	-
POOL2	-	-
CONV3-512	-	-
ReLU	-	-
CONV3-512	-	-
ReLU	-	-
POOL2	-	-
CONV3-512	-	-
ReLU	-	-
CONV3-512	-	-
ReLU	-	-
POOL2	-	-
FC-4096	-	-
FC-4096	-	-
FC-1000	-	-
SOFTMAX	-	-

2. What is the total number of parameters? (use a calculator for this one)
3. What percentage of the weights are found in the fully-connected layers?



Part 2 - Code Assignments

- You must write your code in this notebook and save it with the output of all of the code cells.
- Additional text can be added in Markdown cells.
- You can use any other IDE you like (PyCharm, VSCode...) to write/debug your code, but for the submission you must copy it to this notebook, run the code and save the notebook with the output.

Tips

1. Uniformly distributed tensors - `torch.Tensor(dim1, dim2, ..., dimN).uniform_(-1, 1)`

2. Separation to **validation set** in PyTorch - [See example here](#).

```
In [1]: # imports for the practice (you can add more if you need)
import os
import numpy as np
import pandas as pd
import torch
import torch.nn as nn
from torch.utils.data import TensorDataset, DataLoader
import torchvision
import matplotlib.pyplot as plt
from sklearn import preprocessing
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
# %matplotlib ipynb
%matplotlib inline

seed = 211
np.random.seed(seed)
torch.manual_seed(seed)
```

Out[1]: <torch._C.Generator at 0x7d0600973450>



Task 1 - The Importance of Activation and Initialization

In this task, we are going to use $x \in \mathcal{R}^{512}$ and simple neural network that outputs $f(x) \in \mathcal{R}^{512}$. The network will have 100 layers with 512 units in each layer.

1. We initialize the weights from a unit normal distribution. Run the following code cell and explain what happens. Add a short piece of code that locates when it happens (hint: use `torch.isnan()`). **Print** the layer number.
2. We can demonstrate that at a given layer, the matrix product of inputs x and weight matrix a that is initialized from a standard normal distribution will, on average, have a standard deviation very close to the square root of the number of input connections. For our example, with 512 dimensions, show that for 10,000 multiplications of a and x , the empirical standard deviation is similar to the square root of the number of input connections. Use the unbiased version:

$$\hat{std} = \sqrt{\frac{\sum_{i=1}^{10000} \frac{1}{N} \sum_{j=1}^N y^2}{10000}},$$

where $y = ax$ and N is the number of input connections. **Print** the mean, std and the square root of the number of input connections.

3. For the code from 1, normalize the weight initialization by the square root of the input connections. How does that change the outcome? **Print** the mean and std after the modification.
4. Add a `tanh()` activation after each layer for the code from 1. **Print** the mean and std after the modification. Explain the result.

5. Xavier initialization sets a layer's weights to values chosen from a random uniform distribution that's bounded between

$$\pm \sqrt{\frac{6}{n_i + n_{i+1}}}$$

where n_i is the number of incoming network connections, or "fan-in," to the layer, and n_{i+1} is the number of outgoing network connections from that layer, also known as the "fan-out". Glorot and Bengio believed that Xavier weight initialization would maintain the variance of activations and back-propagated gradients all the way up or down the layers of a network and demonstrated that networks initialized with Xavier achieved substantially quicker convergence and higher accuracy.

Implement **Xavier Uniform** as `xavier_init(fan_in, fan_out)`, a function that returns a tensor initialized according to **Xavier Uniform**. Use it on the simple network from 1 with `tanh` activation. **Print** the mean and std after the modification.

6. If you try to replace the `tanh` activation with `relu` activation in section 5, you will see very different results. Xavier strives to achieve activation outputs of each layer to have a mean of 0 and a standard deviation around 1, on average. When using a ReLU activation, a single layer will, on average have standard deviation that's very close to the square root of the number of input connections, **divided by the square root of two** ($\sqrt{\frac{512}{2}}$ in our example). **Kaiming He et. al.** proposed an initialization scheme that's tailored for deep neural nets that use these kinds of asymmetric, non-linear activations. Implement **Kaiming Normal** as `kaiming_init(fan_in, fan_out)`, a function that returns a tensor initialized according to **Kaiming Normal** (use `fan_in` mode). Use it on the simple network from 1 with `relu` activation. **Print** the mean and std after the modification. What happens when you use Xavier with ReLU activation?

```
In [ ]: x = torch.randn(512)
        for i in range(100):
            a = torch.randn(512, 512)
            x = a @ x
        print(x.mean(), x.std())
```

tensor(nan) tensor(nan)

```
In [ ]: """
        Your Code Here - Use as many blocks as you need
        """
```



Task 2 - MLP-based Deep Classifier

In this task you are going to design and train your first neural network for classification.

For this task, we will use the "[MAGIC Gamma Telescope Data Set](#)". Cherenkov gamma telescope observes high energy gamma rays, taking advantage of the radiation emitted

by charged particles produced inside the electromagnetic showers initiated by the gammas, and developing in the atmosphere. This Cherenkov radiation (of visible to UV wavelengths) leaks through the atmosphere and gets recorded in the detector, allowing reconstruction of the shower parameters. The available information consists of pulses left by the incoming Cherenkov photons on the photomultiplier tubes, arranged in a plane, the camera.

Depending on the energy of the primary gamma, a total of few hundreds to some 10000 Cherenkov photons get collected, in patterns (called the shower image), allowing to discriminate statistically those caused by primary gammas (**signal**) from the images of hadronic showers initiated by cosmic rays in the upper atmosphere (**background**).

Our data has 10 features and 2 classes (signal and background).

1. Load the MAGIC dataset stored in `magic04.data` and display the first 5 features (just run the cell).
2. Separate the data to train, validation and test, reserve 10% of the data for validation and 20% for test.
3. Perform pre-processing steps of your choice and convert the class label from `str` to `int` (for example, `y_train = np.array([0 if y_train[i] == 'g' else 1 for i in range(len(y_train))]).astype(np.int)`).
4. Train a Logistic Regression model from `sklearn` as a baseline for our neural network (only for this section use both the train and validation sets for training the classifier). **Print the test accuracy.**
5. Convert the `numpy` arrays to `torch` tensors with `TensorDataset` as done in the tutorial.
6. Design a **MLP** to classify the data. Optimize the hyper-parameters of your model using the accuracy on the validation set, and when you are satisfied with the model train it on both the train and validation sets and evaluate it on the test set. **You need to reach at least 85% accuracy on the test set, and 87% for a full grade.**
 - You have a free choice of architecture, optimizer, learning scheduler, initialization, regularization and activations.
 - The loss criterion is binary cross entropy: `nn.BCEWithLogitsLoss()` (performs `sigmoid` for you) or `nn.BCELoss` (you need to apply `sigmoid` on the network output yourself).
 - In a Markdown block, write down the chosen architectures and all the hyper-parameters.
 - Make sure to describe any design choice that you used to improve the performance (e.g. if you used a certain regularization or layer, mention it and describe why you think it helped).
 - **Plot** the loss curves (and any other statistic you want) as a function of epochs/iterations. **Print** the final performance.
 - **Print** the test accuracy.
7. Pick **2** initializations of your choosing and change the initialization of the linear layers and re-train the model (with the same optimal hyper-parameters you found). You can pick an initialization of your choosing from :

<https://pytorch.org/docs/stable/nn.init.html> . See example below how to use. **Print** the change in accuracy for both changes (you should end up with 3 results - original, `init 1` and `init 2`).

```
In [ ]: # loading the data
col_names = ['fLength', 'fWidth', 'fSize', 'fConc', 'fConc1', 'fAsym', 'fM3Long']
feature_names = ['fLength', 'fWidth', 'fSize', 'fConc', 'fConc1', 'fAsym', 'fM3Long']
data = pd.read_csv("./magic04.data", names=col_names)
X = data[feature_names]
Y = data['class']
data.head()
```

```
In [ ]: # separate to train, test
"""
Your Code Here
"""
```

```
In [ ]: # pre-processing and converting labels to integers
"""
Your Code Here
"""
```

```
In [ ]: # training a Logistic Regression baseline - complete the code with your variable
logistic_model = LogisticRegression(solver='lbfgs')
y_pred = logistic_model.fit(X_train_prep, y_train_np).predict(X_train_prep)
print("Number of mislabeled points %d out of %d total points." % ((y_train_np != y_pred).sum(), y_train_np.size))
print("Logistic Regression Model accuracy = ", logistic_model.score(X_test_prep, y_test_np))
```

```
In [ ]: # create TensorDataset from numpy arrays
"""
Your Code Here
"""
```

```
In [ ]: # model, hyper-parameters and training
"""
Your Code Here - add as many blocks as you wish
"""
```

```
In [ ]: # example of weight initialization
import torch.nn as nn
class MyModel(nn.Module):
    def __init__(self, parameters):
        super(MyModel, self).__init__()
        # model definitions/blocks
        # ...
        # custom initialization
        self.init_weights()

    def init_weights(self):
        for m in self.modules():
            if isinstance(m, nn.Linear):
                # pick initialization: https://pytorch.org/docs/stable/nn.init.html#examples
                # nn.init.kaiming_normal_(m.weight, mode='fan_out', nonlinearity='relu')
                # nn.init.kaiming_normal_(m.weight, mode='fan_in', nonlinearity='relu')
                # nn.init.normal_(m.weight, 0, 0.005)
                # don't forget the bias term (m.bias)
```

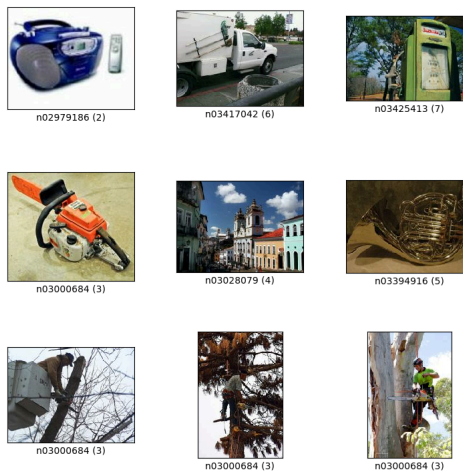
```
def forward(self, x):
    # ops on x
    # ...
    # output = f(x)
    return output
```



Task 3 - Design a CNN

In this task you are going to design a deep convolutional neural network to classify 10 classes from Imagenet (tench, English springer, cassette player, chain saw, church, French horn, garbage truck, gas pump, golf ball, parachute) - **The Imagenette Dataset**.

- 10 classes, 1 for each object.
- 9469 images for training and 3925 for testing (70/30 separation).
- We will use a downscaled version where the images are resized to 64×64 resolution.



1. Load the the Imagenette dataset with PyTorch using

```
torchvision.datasets.Imagenette( root='./datasets', split='train',
size='160px', download=True, transform=transform_train) , where
split is either 'train' or 'val' , you can read more here:
```

<https://pytorch.org/vision/main/generated/torchvision.datasets.Imagenette.html#torchvi>

. Use the `transform` parameter to resize the images to 64×64 (for train, validation and test) and convert the data to tensors, e.g.,

```
transform_test=transforms.Compose([
    transforms.Resize((64, 64)),
    transforms.ToTensor(),])
```

Display 5 images from the train set.

Train, Validation and Test Split for torchvision Datasets

2. Design a Convolutional Neural Network (CNN) to classify classes from the images.

- You are **not allowed** to use `BatchNorm` in your architecture, but can use any other normalization (`GroupNorm` , `LayerNorm` , and etc..).
- Describe the chosen architecture, how many layers? What activations did you choose? What are the filter sizes? Did you use fully-connected layers (if you did, explain their sizes)?
- What is the input dimension? What is the output dimension?
- Calculate the number of parameters (weights) in the network. What is the model size in MegaBytes (MB)? (see the convolution tutorial). **Print** these numbers.

3. Train the classifier (preferably on a GPU - use Colab for this part if you don't have a GPU).

- **DO NOT USE ANY IMAGE AUGMENTATIONS IN THIS PART** (You can still use `Normalize` if you wish, but no cropping, flipping and etc...).
- You are not allowed to use pre-trained models (i.e., no transfer learning, only learning from scratch).
- Describe the hyper-parameters of the model (batch size, epochs, optimizer, learning rate, scheduler....). How did you tune your model? Did you use a validation set to tune the model?
- What is the final accuracy on the test set? **Print** it.
 - You need to reach at least 73% accuracy in this section, and 78% for maximum points in section 5.
- **Plot** the loss curves (and any other statistic you calculate) as a function of epochs/iterations.

4. For the trained classifier, what is the accuracy on the test set when each test image is added a small noise $a = (0.05, 0.01, 0.005)$:

$$\text{image} + a \times \mathcal{N}(0, 1).$$

Print the result for each value of a .

5. Retrain the classifier, but this time use data augmentations of your choosing. Briefly explain what augmentation you chose and how it works. Did the test accuracy improve? **Print** the result.

- You can use transformations available in `torchvision.transforms` as shown in the tutorial.
- You are welcome to use `kornia` for the augmentations (**2 points bonus**, maximal grade is still 100).
- **Plot** the loss curves (and any other statistic you want) as a function of epochs/iterations.

```
In [ ]: """
Your Code Here - add as many blocks as you wish
"""
```



Credits

- Icons made by [Becris](#) from www.flaticon.com
- Icons from [Icons8.com](https://icons8.com) - <https://icons8.com>
- Datasets from [Kaggle](https://www.kaggle.com/) - <https://www.kaggle.com/>