# ECE 046211 - Technion - Deep Learning

## HW1 - Optimization and Automatic Differentiation

### 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 garde: 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_hw1_id1_id2.ipynb`.
  - If you answered the questionss in a different file you should submit a `.zip` file with the name `ece046211_hw1_id1_id2.zip` with content:
    - `ece046211_hw1_id1_id2.ipynb` - the code tasks
    - `ece046211_hw1_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](Anaconda) or **online** on [Google Colab](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` $\rightarrow$ `Change Runtime Type` $\rightarrow$ `GPU`.
  2. Python IDE such as [PyCharm](PyCharm) or [Visual Studio Code](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](https://github.com/taldatech/ee046211-deep-learning)) to help you get everything installed.
- If you need any technical assistance, please go to our forum and describe your problem (preferably with images).

## Agenda

- [Part 1 - Theory](Part 1 - Theory)

# Part 1 - Theory

---

- You can choose whether to answser 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.

- $\LaTeX$ Cheat-Sheet (to write equations)

    - Another Cheat-Sheet

# Question 1 - Convergence of Gradient Descent

---

Recall from the lecture notes:

- **Definition**: A function $f$ is $\beta$-smooth if:
$$\forall w_1, w_2 \in \mathbb{R}^d : ||\nabla f(w_1) - \nabla f(w_2)|| \leq \beta ||w_1 - w_2||$$
- **Lemma**: If $f$ is $\beta$-smooth then
$$f(w_1) - f(w_2) - \nabla f(w_2)^T (w_1 - w_2) \leq \frac{\beta}{2} ||w_1 - w_2||^2$$

Prove the lemma.

Hints:

- Represent $f$ as an integral: $f(x) - f(y) = \int_0^1 \nabla f(y + t(x - y))^T (x - y) dt$
- Make use of Cauchy-Schwarz.

## Question 2 - Optimization and Gradient Descent

We consider the **SGD (stochastic gradient descent)** model we saw in class, for $(w \in \mathbb{R})$, with the quadratic cost function:

$$f(w) = \frac{1}{2} \sum_{n=1}^{N} h_n w^2$$

At each iteration $(t)$ of SGD, we randomly sample an index $(n(t) \in 1, \ldots, N)$, and update according to

$$w(t) = w(t-1) - \eta h_{n(t)} w(t-1)$$

where the random sampling is independent at each step. Here $(\eta > 0)$ is the learning rate. In addition, define

$$h \triangleq \frac{1}{N} \sum_{n=1}^{N} h_n, \quad \rho \triangleq \frac{1}{N} \sum_{n=1}^{N} h_n^2 - h^2$$

1. Write the update equation for $\mathbb{E}[w(t)]$ (that is, $\mathbb{E}[w(t)]$ as a function of $\mathbb{E}[w(t-1)]$, $h$, and $\eta$).

2. Find a condition on $\eta$ (as a function of $h$) such that $\mathbb{E}[w(t)]$ converges to the minimum at the maximal rate.

3. What is the range of $\eta$ values (an inequality as a function of $h$) for which $\mathbb{E}[w(t)]$ necessarily converges to the minimum?

4. Write the update equation for $\mathbb{E}[w^2(t)]$ (that is, $\mathbb{E}[w^2(t)]$ as a function of $\mathbb{E}[w^2(t-1)]$, $h$, $\rho$, and $\eta$).

5. Find a condition on $\eta$ (as a function of $h$ and $\rho$) such that $\mathbb{E}[w^2(t)]$ converges to the minimum at the maximal rate.

6. What is the range of $\eta$ values (an inequality as a function of $h$ and $\rho$) for which $\mathbb{E}[w^2(t)]$ necessarily converges to the minimum?

7. What is the range of $\eta$ values (an inequality as a function of $h$ and $\rho$) for which $\mathbb{E}[f(w(t))]$ necessarily converges to the minimum?

8. Write the update equation for $\log(w(t))$, and use it to prove that, with probability 1,

$$\lim_{t\to\infty} \frac{1}{t}\log(w(t)) = q(\eta) \triangleq \frac{1}{N}\sum_{n=1}^{N}\log(1 - \eta h_n)$$

*Hint:* Use the Law of Large Numbers, i.e., if we have i.i.d. samples $X_n * n = 1^N$ drawn from a distribution with mean $\mathbb{E}[X]$ and bounded variance, then

$$\lim *N \to \infty \frac{1}{N}\sum_{n=1}^{N} X_n = \mathbb{E}[X]$$

with probability 1.

9. What is the condition that must hold on $q(\eta)$ such that $w(t)$ converges to 0 with probability 1?

# Question 3 - Efficient Differentiation

We wish to optimize a loss function $\mathcal{L}(\mathbf{w})$ for $\mathbf{w} \in \mathbb{R}^d$ using Gradient Descent (GD) with some step size schedule $\eta_t$

$$(1) \quad \forall t = 1, 2, \dots : \mathbf{w}(t) = \mathbf{w}(t-1) - \eta_t \nabla \mathcal{L}(\mathbf{w}(t-1))$$

initialized from some $\mathbf{w}(0)$. We would like to learn the best step size schedule using GD. **Hint:** throughout this question, you should use the *chain rule*.

1. Suppose we can consider each $\eta_t$ as a separate parameter for each $t$. We initialize this parameter with $\eta_0$ and update $\eta_{t-1}$ with a GD step on $\mathcal{L}(\mathbf{w}(t-1))$

$$(2) \quad \eta_t = \eta_{t-1} - \alpha_t \frac{\partial \mathcal{L}(\mathbf{w}(t-1))}{\partial \eta_{t-1}}$$

for every step of eq. $(1)$, where $\alpha_t$ is the another step size. Calculate $\partial \mathcal{L}(\mathbf{w}(t-1))/\partial \eta_{t-1}$ as a function of the loss gradients $\nabla \mathcal{L}(\mathbf{w}(t-1))$ and $\nabla \mathcal{L}(\mathbf{w}(t-2))$.

2. Now suppose we want to similarly update $\alpha_{t-1}$ using GD step on $\mathcal{L}(\mathbf{w}(t-1))$ every step of eq. $(2)$ with update step $\kappa_t$

$$\alpha_t = \alpha_{t-1} - \kappa_t \frac{\partial \mathcal{L}(\mathbf{w}(t-1))}{\partial \alpha_{t-1}}.$$

Calculate $\partial \mathcal{L}(\mathbf{w}(t-1))/\partial \alpha_{t-1}$ as a function of $\{\nabla \mathcal{L}(\mathbf{w}(t-k))\}_{k=1}^{3}$.

3. Now we wish to update $(\eta_{t-1}, \eta_{t-2})$ by doing a GD step on $\mathcal{L}(\mathbf{w}(t-1))$

$$(3) \quad (\eta_{t+1}, \eta_t) = (\eta_{t-1}, \eta_{t-2}) - \alpha_t \left(\frac{\partial \mathcal{L}(\mathbf{w}(t-1))}{\partial \eta_{t-1}}, \frac{\partial \mathcal{L}(\mathbf{w}(t-1))}{\partial \eta_{t-2}}\right)$$

every two steps of eq. $(1)$. Calculate the derivative $\frac{\partial \mathcal{L}(\mathbf{w}(t-1))}{\partial \eta_{t-2}}$ as a function of $\eta_{t-1}$, $\{\nabla \mathcal{L}(\mathbf{w}(t-k))\}_{k=1}^{3}$, and $\nabla^2 \mathcal{L}(\mathbf{w}(t-2))$.

4. Now we wish again to update $(\eta_t, \eta_{t+1}, \ldots, \eta_{t+T})$ by doing a GD step on $\mathcal{L}(\mathbf{w}(t-1))$ every $T$ steps of eq. $(1)$

$$(4) \quad (\eta_{t+T}, \ldots, \eta_t) = (\eta_{t-1}, \ldots, \eta_{t-1-T}) - \alpha_t \left( \frac{\partial \mathcal{L}(\mathbf{w}(t-1))}{\partial \eta_{t-1}}, \ldots, \frac{\partial \mathcal{L}(\mathbf{w}(t-1))}{\partial \eta_{t-1-T}} \right)$$

Calculate the derivative $\frac{\partial \mathcal{L}(\mathbf{w}(t-1))}{\partial \eta_{t-\tau}}$ as a function of $\{\eta_{t-k}, \nabla^2 \mathcal{L}(\mathbf{w}(t-k-1))\}_{k=1}^{\tau-1}$, $\nabla \mathcal{L}(\mathbf{w}(t-1))$ and $\nabla \mathcal{L}(\mathbf{w}(t-\tau-1))$.

5. Compare this approach (eq. $(4)$) with $T > 1$) to the first one (eq. $(2)$). Name one advantage for each approach. Hints: Think of compuptional complexity, ease of optimization, suitability of the objective.

## Question 4 - Automatic Differentiation

Consider the scalar function:

$$f = \exp(\exp(x) + \exp(x)^2) + \sin(\exp(x) + \exp(x)^2)$$

1. Write down the derivative w.r.t. $x$ explicitly, i.e., $\frac{df}{dx}$

2. We define the following intermediate variables:

$$a = \exp(x)$$
$$b = a^2$$
$$c = a + b$$
$$d = \exp(c)$$
$$e = \sin(c)$$
$$f = d + e$$

Draw a graph picturing the relationship between all variables (called the **computation graph**).

3. Using the graph, write down the derivatives of the individual terms, working backwards to compute the derivative of $f$ (i.e., write down the derivatives $\frac{df}{dd}, \frac{df}{de}, \ldots, \frac{df}{dx}$)

## Question 5 - Automatic Differentiation 2

Write down the chain rule in the dual numbers representation for the following:

$$f(g(h(x + \epsilon x')))$$

What is $\frac{df(x)}{dx}$?

## Part 2 - Code Assignments

- You must write your code in this notebook and save it with the output of aall 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.

```
 1 # imports for the practice (you can add more if you need)
 2 import os
 3 import numpy as np
 4 import pandas as pd
 5 import torch
 6 import matplotlib.pyplot as plt
 7 from mpl_toolkits.mplot3d import Axes3D
 8 from matplotlib.colors import LogNorm
 9 from sklearn.datasets import load_iris
10 seed = 211
11 np.random.seed(seed)
12 torch.manual_seed(seed)
13 # %matplotlib notebook
14 %matplotlib inline
```

## Task 1 - The Beale Function

The Beale function is defined as follows:

$$f(x, y) = (1.5 - x + xy)^2 + (2.25 - x + xy^2)^2 + (2.625 - x + xy^3)^2$$

1. What is the global minima of this function?
2. Implement the Beale function: `beale_f(x,y)`.
3. Implement a function, `beale_grads(x,y)` that returns the gradients of the Beale function.
4. 3D plot the Beale function wit the global minima you found. Use Matplotlib's `ax.plot_surface(x_mesh, y_mesh, z, norm=LogNorm(), rstride=1, cstride=1, edgecolor='none', alpha=.8, cmap=plt.cm.jet)` for the function, and `ax.plot(x, y, f(x, y), 'r*', markersize=20)` for the minima.

5. 2D plot the contours with `ax.contour(x_mesh, y_mesh, z, levels=np.logspace(-.5, 5, 35), norm=LogNorm(), cmap=plt.cm.jet)` and the minima with `ax.plot(x, y, 'r*', markersize=20)`.

Your Answers Here

```
 1 # Set the manually calculated minima
 2 min_x = None
 3 min_y = None
 4
 5 def beale_f(x, y):
 6     value = None
 7     """
 8     Your Code Here
 9     """
10     return value
11
12 def beale_grads(x, y):
13     dx, dy = None, None
14     """
15     Your Code Here
16     """
17
18     grads = np.array([dx, dy])
19     return grads
```

```
1 minima = np.array([min_x, min_y])
2 beale_res = beale_f(*minima)
3 grads_res = beale_grads(*minima)
4 print(f"minima (1x2 row vector shape): {minima}")
5 print(f"beale_f output: {beale_res}")
6 print(f"beale_grad output: {grads_res}")
```

## `</>` Task 2 - Building an Optimizer - Adam

In this task, you are going to implement the Adam optimizer. We are giving the skeleton of the code and the description of the methods, and you need to implement the optimizer.

Recall the Adam update rule:
$$m_{k+1} = \beta_1 m_k + (1 - \beta_1)\nabla f(w^k) = \beta_1 m_k + (1 - \beta_1)g_k$$
$$v_{k+1} = \beta_2 v_k + (1 - \beta_2)(\nabla f(w^k))^2 = \beta_2 v_k + (1 - \beta_2)g_k^2$$

Then, they use an **unbiased** estimation:
$$\hat{m}_{k+1} = \frac{m_{k+1}}{1 - \beta_1^{k+1}}$$

$$\hat{v}_{k+1} = \frac{v_{k+1}}{1 - \beta_2^{k+1}}$$

(the $\beta$'s are taken with the power of the current iteration)

$$w_{k+1} = w_k - \frac{\alpha}{\sqrt{\hat{v}_{k+1}} + \epsilon} \hat{m}_{k+1}$$

- $\epsilon$ deafult's is $10^{-8}$

1. Implement `class AdamOptimizer()`.

   - `function` is the Python function you want to optimize.
   - `gradients` is the Python function that returns the gradients of `function`.
   - `x_init` and `y_init` are the initialization points for the optimizer.
   - Save the `path` of the optimizer (the minima points the optimizer visits during the optimization).
   - Stopping criterion: change in minima `<1e-7`.
   - **You can change the class however you wish, you can remove/add variables and methods as you wish**

2. For `x_init=0.7, y_init=1.4, learning_rate=0.1, beta1=0.9, beta2=0.999`, optimize the Beale function. Plot the results **with the path taken** (better do it on the 2D contour plot).

3. Choose different initialization and learning rate and show the results as in 2.

```
1 class AdamOptimizer():
2     def __init__(self, function, gradients, x_init=None, y_init=None,
3                  learning_rate=0.001, beta1=0.9, beta2=0.999, epsilon=1e-8):
4         self.f = function
5         self.g = gradients
6         scale = 3.0
7         self.current_val = np.zeros([2])
8         if x_init is not None:
9             self.current_val[0] = x_init
10        else:
11            self.current_val[0] = np.random.uniform(low=-scale, high=scale)
12        if y_init is not None:
13            self.current_val[1] = y_init
14        else:
15            self.current_val[1] = np.random.uniform(low=-scale, high=scale)
16        print("x_init: {:.3f}".format(self.current_val[0]))
17        print("y_init: {:.3f}".format(self.current_val[1]))
18
19        self.lr = learning_rate
20        self.grads_first_moment = np.zeros([2])
21        self.grads_second_moment = np.zeros([2])
22        self.beta1 = beta1
23        self.beta2 = beta2
24        self.epsilon = epsilon
25
```

```python
26          # for accumulation of loss and path (w, b)
27          self.z_history = []
28          self.x_history = []
29          self.y_history = []
30
31
32      def func(self, variables):
33          """Beale function.
34
35          Args:
36            variables: input data, shape: 1-rank Tensor (vector) np.array
37              x: x-dimension of inputs
38              y: y-dimension of inputs
39
40          Returns:
41            z: Beale function value at (x, y)
42          """
43
44      def gradients(self, variables):
45          """Gradient of Beale function.
46
47          Args:
48            variables: input data, shape: 1-rank Tensor (vector) np.array
49              x: x-dimension of inputs
50              y: y-dimension of inputs
51
52          Returns:
53            grads: [dx, dy], shape: 1-rank Tensor (vector) np.array
54              dx: gradient of Beale function with respect to x-dimension of inputs
55              dy: gradient of Beale function with respect to y-dimension of inputs
56          """
57
58
59      def weights_update(self, grads, time):
60          """Weights update using Adam.
61
62            g1 = beta1 * g1 + (1 - beta1) * grads
63            g2 = beta2 * g2 + (1 - beta2) * grads ** 2
64            g1_unbiased = g1 / (1 - beta1**time)
65            g2_unbiased = g2 / (1 - beta2**time)
66            w = w - lr * g1_unbiased / (sqrt(g2_unbiased) + epsilon)
67          """
68
69      def history_update(self, z, x, y):
70          """Accumulate all interesting variables
71          """
72
73
74      def train(self, max_steps):
75
```

```
1 """
2 Your Code Here
3 """
```

```
1 opt = AdamOptimizer(beale_f, beale_grads, x_init=0.7, y_init=1.4, learning_rate=0.1,
```

```
1 %time
2 opt.train(1000)
3 print("Global minima")
4 print("x*: {:.2f}  y*: {:.2f}".format(minima[0], minima[1]))
5 print("Solution using the gradient descent")
6 print("x: {:.4f}  y: {:.4f}".format(opt.x, opt.y))
```

```
1 # plot the Beale function values during the optimization
```

```
1 # plot the optimization path
2 path = opt.path
```

## Task 3 - PyTorch Autograd

For the function from the theory practice:
$$f = \exp(\exp(x) + \exp(x)^2) + \sin(\exp(x) + \exp(x)^2)$$

1. Implement it and its dervative (explicitly) using `torch`.
2. Define a scalar tensor `x` and use `autograd` to calculate the derivative w.r.t $x$. Does the result correspond to the output of the function the calculates the derivative explicitly?

```
1 def f(x):
2     f_val = None
3     """
4     Your Code Here
5     """
6     return f_val
7
8 def derv_f(x):
9     derv_val = None
10    """
11    Your Code Here
12    """
13    return derv_val
```

```
1 x = torch.tensor(0.5, requires_grad=True)
2 print(x)
3 f_res = f(x)
4 f_manual_grad = derv_f(x.detach())
```

```
 5
 6 """
 7 Your Code Here
 8 """
 9 # Calculate with torch autograd
10 f_autograd = None
11
12
13 print(f_manual_grad)
14 print(f_autograd)
```

</> **Task 4 - Low Rank Matrix Factorization**

Consider the following optimization problem:

$$\min_{\hat{U},\hat{V}} ||A - \hat{U}\hat{V}||_F^2$$

Where $A \in \mathcal{R}^{m \times n}, \hat{U} \in \mathcal{R}^{m \times r}, \hat{V} \in \mathcal{R}^{r \times n}$ and $r < min(m, n)$ ($r$ is the rank of the matrix). $|| \cdot ||_F^2$ denotes the Frobenius norm.

1. Implement a function, `gd_factorize_ad(A, rank, num_epochs=1000, lr=0.01)`, that given a 2D tensor `A` and a `rank`, will calculate the low-rank factorization of `A` using **gradient decsent**. Compute and apply all the gradients of $\hat{U}$ and of $\hat{V}$ once per epoch. $\hat{U}$ and $\hat{V}$ should be initially created with uniform random values. Use PyTorch's `autograd` for the gradients.

   ○ To compute the squared Frobenius norm loss (reconstruction loss), use `torch.nn.functional.mse loss with reduction='sum'`.

2. Use the provided `data` of the Iris dataset of 150 instances and 4 features. Apply `gd_factorize_ad` to compute the 2-rank matrix factorization of `data`. What is the reconstruction loss?

```
1 df = load_iris(as_frame=True).data # option 1
2 # df = pd.read_csv('./iris.data', header=None) # option 2
3 data = torch.tensor(df.iloc[:, [0, 1, 2, 3]].values)
4 data = data - data.mean(dim=0)
```

```
1 def gd_factorize_ad(A, rank, num_epochs=1000, lr=0.01):
2     # initialize
3     U = None
4     V = None
5
6     """
7     Your Code Here
8     """
```

```
 9
10     # implement gradient descent
11     for epoch in range(num_epochs):

12

13         """
14         Your Code Here
15         """

16         loss = None
17         if epoch % 5 == 0:
18             print(f'epoch: {epoch}, loss: {loss}')
19     return U, V
```

```
1 U, V = gd_factorize_ad(data.float(), rank=2, num_epochs=1000, lr=0.01)
```

# Credits

- Icons made by Becris from www.flaticon.com
- Icons from Icons8.com - https://icons8.com
- Datasets from Kaggle - https://www.kaggle.com/