

dpoints

**ECE 364: Programming Methods for Machine Learning,
Spring 2025**

Midterm 1 – March 11, 2025

-
- **You will have 75 minutes (1.25 hours) to solve all the problems. Most have multiple parts.** Don't spend too much time on questions you don't understand and focus on answering as much as you can!
 - ***BUDGET YOUR TIME WISELY***. I highly recommend working on the questions you know first and the questions you need to think about second.
 - *No* resources are allowed for use during the exam except a cheatsheet and scratch paper on the back of the exam. **Do not tear out the cheatsheet or the scratch paper!** It messes with the auto-scanner.
 - You should write your answers *completely* in the space given for the question. We will not grade parts of any answer written outside of the designated space.
 - Please *use a dark-colored pen* unless you are *absolutely* sure your pencil writing is forceful enough to be legible when scanned. We reserve the right to take off points if we have difficulty reading the uploaded document.
 - **Don't cheat.** C'mon, be cool, be honest.
 - **Good luck!**
-

Name: _____

NetID: _____

Date: _____

1. Lost in the Layers: Navigating Data Slices

(24 points)

For each of the code segments, answer the following questions based on the final state of the variables.

```
(a) import torch
2 a = torch.tensor([[0,1,2],[3,4,5],[6,7,8],[9,10,11]])
3 b = torch.arange(3).view(1,3)
4 c = a*b
5 c += 1
```

i. What will be the shape of `c`?

Solution:

4×3 (4 rows by 3 columns)

ii. What does `b` contain (note dimension in addition to value)?

Solution:

`[[0, 1, 2]]`

```
(b) import torch
2 d = torch.arange(100).view(10,10) #arange(N) gives int vector [0...N-1]
3 e = d[0:2, 1:6:4]
4 f = e[0,:]
5 g = f.add_(1)[:2]
6 h = g.t() @ g
```

i. What is the value of `h`?

Solution:

40

ii. What does `e` contain?

Solution:

$$\begin{bmatrix} 2 & 6 \\ 11 & 15 \end{bmatrix}$$

2. Feeling out of place

(5 points)

Suppose we want to find the value of x where $\ln(x^2) = 3$. We write the following gradient descent code to find the issue:

```
1 x_gd = torch.tensor(6.0, requires_grad=True)
2 target=3
3 alpha = 0.001
4 epochs = 20000
5 for i in range(epochs):
6     f = torch.log(torch.pow(x_gd,2))
7     loss = (target-f)**2
8     loss.backward()
9     with torch.no_grad():
10        x_gd = x_gd - alpha*x_gd.grad
11        x_gd.grad = None
```

...but there's an error! After much debugging, you narrow the issue down to the highlighted line above. What do you need to change this line to, so that the gradient descent code can work appropriately?

Solution:

The issue is that you are not doing an in-place operation, so after the first iteration, you are basically just generating a new tensor which defaults to `requires_grad=False` and hence, the next time you try to calculate `loss.backward()`, an error gets thrown because there are no gradients to compute. You need to change the line to `x_gd -= alpha*w_gd` which is an in-place operation and thus the tensor is preserved with the options you declared previously. If you want, you can play around with the code below:

```
1
2 import torch
3 import numpy as np
4 import matplotlib.pyplot as plt
5
6 x_vals = []
7 f_vals = []
8 x_gd = torch.tensor(6.0, requires_grad=True)
9 target=3
10 alpha = 0.001
11 epochs = 20000
12 for i in range(epochs):
13     f = torch.log(torch.pow(x_gd,2))
14     loss = (target-f)**2
15     loss.backward()
16     with torch.no_grad():
17         x_gd -= alpha*x_gd.grad
18         #x_gd = x_gd - alpha*x_gd.grad
19         x_gd.grad = None
20         x_vals.append(x_gd.data.item())
21         f_vals.append(f.data.item())
22
23 fig, ax1 = plt.subplots()
24 color = 'tab:red'
25 ax1.set_xlabel('iterations')
26 ax1.set_ylabel('x values', color=color)
```

```

27 ax1.plot(range(epochs), x_vals, color=color)
28 ax1.tick_params(axis='y', labelcolor=color)
29
30 ax2 = ax1.twinx() # instantiate a second Axes that shares the same x-axis
31 color = 'tab:blue'
32 ax2.set_ylabel('f(x)', color=color) # we already handled the x-label with ax1
33 ax2.plot(range(epochs), f_vals, color=color, linestyle='--')
34 ax2.tick_params(axis='y', labelcolor=color)
35
36 fig.tight_layout() # otherwise the right y-label is slightly clipped
37 plt.show()
38
39 print('x = ' + str(x_gd.data.item()) + ' when f(x) = ' + str(target) )
40

```

3. Jacobians, Hessians, and Why My Brain Hurts (Matrix Calculus) (5 points)

Given:

$$f(x) = x^T x$$

where $x \in \mathbb{R}^n$

Find the solution to the following partial derivative:

$$\frac{\partial f}{\partial x} =$$

You must explain why your answer is correct for full credit.

Solution:

This is right out of your homework and lecture and so I'm just going to copy the solution from there:

$$\frac{\partial x^T x}{\partial x} = 2x$$

Think of it this way, the product of a scalar and its transpose is scalar size:

$$x^T x = x_0^2 + x_1^2 + \dots + x_n^2 = z(x)$$

so measuring the effect each x_i value has on the end sum would yield a vector:

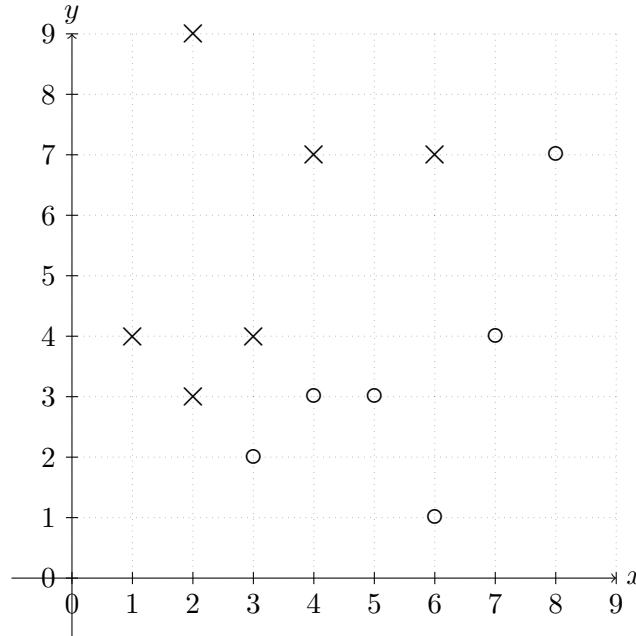
$$\left[\frac{\partial z(x)}{\partial x_0}, \frac{\partial z(x)}{\partial x_1}, \dots, \frac{\partial z(x)}{\partial x_n} \right] = [2x_0, 2x_1, \dots, 2x_n] = 2x$$

4. **The Art of the SVM Margin**

(10 points)

Support vectors are the critical data points in an SVM that lie closest to the decision boundary. They determine both the optimal separating boundary and the margin width between classes. Removing these points from the training set would shift the boundary, potentially altering the classification results.

Consider a SVM that separates the following 2D points into two classes (X and O):



With points:

- **Class X:** $\{(1, 4), (2, 3), (2, 9), (3, 4), (4, 7), (6, 7)\}$
- **Class O:** $\{(3, 2), (4, 3), (5, 3), (6, 1), (7, 4), (8, 7)\}$

- (a) Drawing the Boundary: Sketch the decision boundary for an SVM on the provided graph. Choose the decision boundary that maximizes the margin and minimizes the loss. Additionally, sketch the margin boundaries for each class (these are the lines that pass through the closest data points from each class, positioned half a margin away from the decision boundary). Label each line clearly. You should draw a total of three lines. Include the approximate equation for each line.

Solution:

Class X boundary: line with equation $y = x + 1$. Class O boundary: line with equation $y = x - 1$. Decision boundary: line with equation $y = x$.

- (b) Identifying Support Vectors: For each class (X and O) in the provided data, identify the support vectors and list them below.

Solution:

Class X: $\{(2, 3), (3, 4), (6, 7)\}$, Class O: $\{(3, 2), (4, 3), (8, 7)\}$.

- (c) SVM Loss Function Analysis: Below is the loss function for soft-margin SVM, for a dataset of the form $\mathcal{D} := \{(x_i, y_i)\}_{i=1}^N$ and where $x_i \in \mathbb{R}^n$ and $y_i \in \{+1, -1\}$:

$$\mathcal{L} = \min_{w, \xi} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^n \xi_i$$

where $\xi_i = \max(0, 1 - y_i \cdot Wx_i)^2$. Here $\frac{1}{2} \|w\|_2^2$ is the regularization term, and $C \sum_{i=1}^n \xi_i$ is the slack penalty term. Explain briefly (no more than 4 sentences) what would happen if C is decreased. Include how training and test accuracy might be affected.

Solution:

When $C \rightarrow 0$, the penalty for misclassification becomes negligible compared to the regularization term. Consequently, the SVM focuses on maximizing the margin, even if that means allowing more misclassifications. This results in a wider margin but a decision boundary that may not classify all training points correctly. Overall, the classifier becomes more tolerant of errors, potentially reducing its accuracy on the training data but possibly increasing generalization (test accuracy) if noise is present.

Solution:

A second answer could be that the model simply minimizes the fitting parameters resulting in a flat line with a slope and bias of zero. This would in effect result in a infinite margin which corresponds to the answer above.

5. **Where to Draw the Line? A Classifier's Dilemma**

(10 points)

Let g be the logical **IMPLICATION** function ($A \implies B$), defined on the feature space $\{+1, -1\}^2$, which maps:

- $g(+1, +1) = +1$
- $g(-1, +1) = +1$
- $g(+1, -1) = -1$
- $g(-1, -1) = +1$.

Given a linear classifier $h(x) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$, where $\mathbf{w} \in \mathbb{R}^{1 \times 2}$, $\mathbf{x} \in \mathbb{R}^{2 \times 1}$, and $b \in \mathbb{R}^{1 \times 1}$, give a valid (\mathbf{w}, b) pair that matches the ground truth g . Let $\text{sign}(z) = +1$ for $z \geq 0$ and -1 otherwise. Give your solution and show that it is valid.

Solution:

Let $\mathbf{w} = \begin{bmatrix} -1 & +1 \end{bmatrix}$ and $b = +1$. Then the (\mathbf{w}, b) pair would match the ground truth g for a linear classifier $h(x) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$. We can check by doing the following calculations, seeing that this mapping leads to:

$$h(+1, +1) = \text{sign}\left(\begin{bmatrix} -1 & +1 \end{bmatrix} \cdot \begin{bmatrix} +1 \\ +1 \end{bmatrix} + 1\right) = \text{sign}(+1) = +1 = g(+1, +1)$$

$$h(-1, +1) = \text{sign}\left(\begin{bmatrix} -1 & +1 \end{bmatrix} \cdot \begin{bmatrix} -1 \\ +1 \end{bmatrix} + 1\right) = \text{sign}(+3) = +1 = g(-1, +1)$$

$$h(+1, -1) = \text{sign}\left(\begin{bmatrix} -1 & +1 \end{bmatrix} \cdot \begin{bmatrix} +1 \\ -1 \end{bmatrix} + 1\right) = \text{sign}(-1) = -1 = g(+1, -1)$$

$$h(-1, -1) = \text{sign}\left(\begin{bmatrix} -1 & +1 \end{bmatrix} \cdot \begin{bmatrix} -1 \\ -1 \end{bmatrix} + 1\right) = \text{sign}(+1) = +1 = g(-1, -1)$$

and matches the mapping of g where g is the logical **IMPLICATION** function defined on the feature space $\{+1, -1\}^2$.

6. Stumbling Down the Gradient: My Life Story

(15 points)

Consider the following function

$$f(x, y) = -x \ln x - y \ln y$$

- (a) Determine the gradient $\nabla f(x, y)$.

Solution:

$$\frac{df}{dx} = -1 - \ln x$$
$$\frac{df}{dy} = -1 - \ln y$$

Hence,

$$\nabla f(x, y) = \begin{bmatrix} -1 - \ln x \\ -1 - \ln y \end{bmatrix}$$

- (b) Let the starting point for gradient descent at $k = 0$ be $(x^{(0)}, y^{(0)}) = (1, 1)$ and the step size be $\alpha = e - 1$. Here, e is Euler's number. Apply gradient descent to obtain the values of x and y at iterations $k = 1$ and $k = 2$.

Solution:

Gradient descent update can be written as

$$(x^{(k)}, y^{(k)}) = (x^{(k-1)}, y^{(k-1)}) - \alpha \nabla f(x, y)|_{(x^{(k-1)}, y^{(k-1)})}$$

At $k = 1$, we have

$$\nabla f(x, y)|_{(1,1)} = \begin{bmatrix} -1 & -1 \end{bmatrix}^T$$
$$(x^{(1)}, y^{(1)}) = \begin{bmatrix} 1 \\ 1 \end{bmatrix} - (e - 1) \begin{bmatrix} -1 \\ -1 \end{bmatrix} = \begin{bmatrix} e \\ e \end{bmatrix}$$

At $k = 2$, we have

$$\nabla f(x, y)|_{(e,e)} = \begin{bmatrix} -2 & -2 \end{bmatrix}^T$$
$$(x^{(2)}, y^{(2)}) = \begin{bmatrix} e \\ e \end{bmatrix} - (e - 1) \begin{bmatrix} -2 \\ -2 \end{bmatrix} = \begin{bmatrix} 3e - 2 \\ 3e - 2 \end{bmatrix}$$

7. **torch.nn models, Dataloaders and Optimizers, oh my!** (10 points)

Gradient accumulation is a technique that allows training with larger batch sizes by accumulating gradients over multiple smaller mini-batches before updating model parameters. Normally, in stochastic gradient descent, we compute gradients using `loss.backward()`, update parameters with `optimizer.step()`, and reset gradients using `optimizer.zero_grad()`.

However, if `optimizer.zero_grad()` is not called after each mini-batch, gradients accumulate over multiple iterations. This is useful when a large batch size cannot fit into memory. Instead of processing the full batch at once, we split it into smaller chunks, compute `loss.backward()` for each chunk, and update the model only after accumulating gradients from all chunks.

Now, suppose we want to train with a batch size of 128, but our device can only handle a batch size of 8. Use gradient accumulation to achieve this by completing the code snippet below and correctly placing `loss.backward()`, `optimizer.step()`, and `optimizer.zero_grad()`.

You can assume the following:

- `batch` is a dictionary containing the input and target, where the input has a shape of $B \times M$ and the target has a shape of B . B and M indicate the batch size and number of features.
- The Model (`model`) expects an input of shape $B \times M$ and produces an output of shape B .

```

1 import torch
2 from torch.utils.data import DataLoader
3 from lib.model import Model
4 from lib.data import SimpleDataset
5
6 model = Model()
7 dataset = SimpleDataset()
8 criterion = nn.MSELoss()
9 optimizer = torch.optim.SGD(model.parameters(), lr=0.01, momentum=0.9)
10 # Complete this
11 train_loader =
12
13
14
15 # Complete this
16 num_accumulation_steps =
17
18
19
20 batches_processed = 0
21 optimizer.zero_grad()
22
23 # Add the missing function calls in the below
24 # training code
25 for epoch in range(10):
26     for batch in train_loader: # Can add code between any of these commands
27
28
29         inp, tgt = batch["input"], batch["target"]
30
31
32         output = model(inp)
33
34
35         loss = criterion(output, tgt)
36
37
38
39         batches_processed += 1
40         if batches_processed % num_accumulation_steps == 0:
41             # one more place to fill in code
42
43
44
45
46
47
48     #end of batch loop
49
50
51
52 #end of epoch loop

```

Solution:

```

1 import torch

```

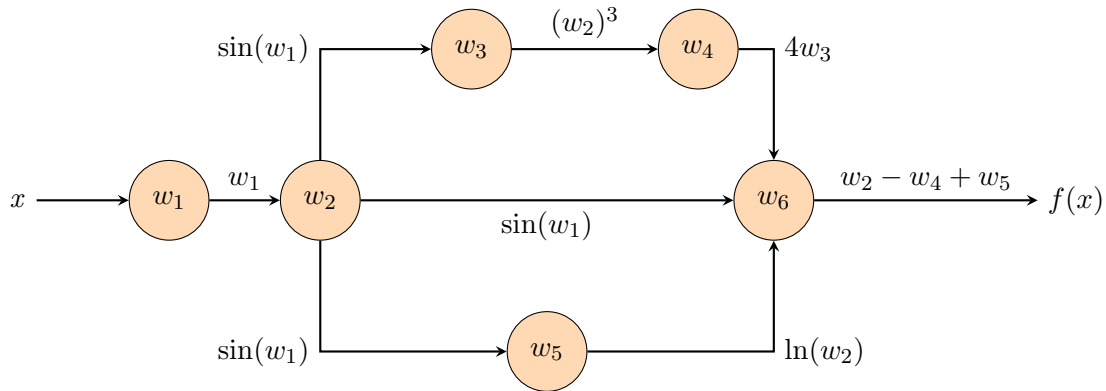
```

2 from lib.model import Model
3 from lib.data import SimpleDataset
4
5 model = Model()
6 dataset = SimpleDataset()
7 train_loader, _ = split_dataset(dataset)
8 criterion = nn.MSELoss()
9 optimizer = torch.optim.SGD(model.parameters(), lr=0.01, momentum=0.9)
10 num_accumulation_steps = 16 # Number of accumulation steps
11
12 batches_processed = 0
13 optimizer.zero_grad()
14 for epoch in range(10):
15     for batch in train_loader:
16         inp, tgt = batch["input"], batch["target"]
17         output = model(inp)
18         loss = criterion(output, tgt)
19         loss.backward()
20         batches_processed += 1
21         if batches_processed % num_accumulation_steps == 0:
22             optimizer.step()
23             optimizer.zero_grad()

```

8. The Computational Graph That Ate My Sanity

(21 points)



(a) Determine the function $f(x)$ represented by the above computational graph.

Solution:

$$\begin{aligned}
 w_2 &= \sin(w_1) = \sin(x) \\
 w_3 &= (w_2)^3 = \sin^3(x) \\
 w_4 &= 4w_3 = 4\sin^3(x) \\
 w_5 &= \ln(w_2) = \ln(\sin(x)) \\
 w_6 &= w_2 - w_4 + w_5 \\
 f(x) &= w_6 = \sin(x) - 4\sin^3(x) + \ln(\sin(x))
 \end{aligned}$$

(b) Determine the partial derivatives of each successor node with respect to its predecessors, e.g., $\frac{\partial w_6}{\partial w_5}$, $\frac{\partial w_6}{\partial w_4}$, $\frac{\partial w_6}{\partial w_2}$, etc.

Solution:

$$\begin{aligned}
 \frac{\partial w_6}{\partial w_5} &= 1 \\
 \frac{\partial w_3}{\partial w_2} &= 3w_2^2 \\
 \frac{\partial w_5}{\partial w_2} &= \frac{1}{w_2} \\
 \frac{\partial w_6}{\partial w_4} &= -1 \\
 \frac{\partial w_2}{\partial w_1} &= \cos(w_1) \\
 \frac{\partial w_4}{\partial w_3} &= 4 \\
 \frac{\partial w_6}{\partial w_2} &= 1 - 12\sin^2(x) + \frac{1}{\sin(x)}
 \end{aligned}$$

(c) Determine the adjoints at each node $\bar{w}_i = \frac{\partial f}{\partial w_i}$.

Solution:

$$\begin{aligned}\bar{w}_6 &= 1 \\ \bar{w}_5 &= \frac{\partial f}{\partial w_5} = \frac{\partial f}{\partial w_6} \frac{\partial w_6}{\partial w_5} = 1 \\ \bar{w}_4 &= \frac{\partial f}{\partial w_4} = \bar{w}_6 \frac{\partial w_6}{\partial w_4} = -1 \\ \bar{w}_3 &= \frac{\partial f}{\partial w_3} = \bar{w}_4 \frac{\partial w_4}{\partial w_3} = -4 \\ \bar{w}_2 &= \frac{\partial f}{\partial w_2} = 1 - 12w_2^2 + \frac{1}{w_2} \\ \bar{w}_1 &= \bar{w}_2 \frac{\partial w_2}{\partial w_1} = \cos(w_1) - 12 \sin^2(w_1) \cos(w_1) + \frac{1}{\sin(w_1)} \cos(w_1)\end{aligned}$$

Solution:

We verified the solutions via PyTorch. Feel free to check!

```

1 import torch
2 import numpy as np
3
4 x = torch.tensor([1.5], requires_grad=True) # make sure gradients are
      computed when backpropagation is called
5 y = torch.tensor([np.pi/3], requires_grad=True)
6
7 w1 = x
8 w2 = torch.sin(w1)
9 w3 = torch.pow(w2, 3)
10 w4 = 4*w3
11 w5 = torch.log(w2)
12 w6 = w2 - w4 + w5
13 f = w6
14
15 # manual gradients
16 with torch.no_grad():
17     # adjoints
18     w6bar = 1
19     w5bar = 1
20     w4bar = -1
21     w3bar = -4
22     w2bar = 1 - (12*(w2**2)) + (1/w2)
23     w1bar = torch.cos(w1) - (12*(torch.sin(w1)**2)*torch.cos(w1)) + (1/
      torch.sin(w1))*torch.cos(w1)
24
25 # automatic gradients via backpropagation
26 w1.retain_grad(), w2.retain_grad(), w3.retain_grad(), w4.retain_grad(), w5
      .retain_grad(), w6.retain_grad() # making sure PyTorch populates all
      gradients
27 f.backward() # initiate backpropagation from f as the seed node
28
29 print("Making sure the overall equation is correct: ")
30 f_manual = torch.sin(x)-4*(torch.sin(x)**3)+torch.log(torch.sin(x))

```

```
31 print('f: Manual = {}, PyTorch = {}'.format(f_manual, f))
32
33 print('Comparing our calculations to PyTorch Autograd:')
34 print('w1: Manual = {}, PyTorch = {}'.format(w1bar, w1.grad))
35 print('w2: Manual = {}, PyTorch = {}'.format(w2bar, w2.grad))
36 print('w3: Manual = {}, PyTorch = {}'.format(w3bar, w3.grad))
37 print('w4: Manual = {}, PyTorch = {}'.format(w4bar, w4.grad))
38 print('w5: Manual = {}, PyTorch = {}'.format(w5bar, w5.grad))
39 print('w6: Manual = {}, PyTorch = {}'.format(w6bar, w6.grad))
40
```

This page is for additional scratch work!

PyTorch Cheatsheet - Part 1

Useful activation function and torch.nn.functional

- Linear function: $y = WX + b$ where W and X are vectors of size N (number of dimensions to the input).

```
torch.nn.Linear(in_features, out_features, bias=True, device=None, dtype=None)
```

- Sigmoid function: $\frac{1}{1+e^{-z}}$ where z is the logit(s).

```
torch.nn.functional.sigmoid(input)
```

- Softmax function: $p(Y = t|x) = \frac{\exp(w_t^T x)}{\sum_{y \in \{0, \dots, C-1\}} \exp(w_y^T x)}$

```
torch.nn.functional.softmax(input, dim=None, _stacklevel=3, dtype=None)[source]
```

Loss Functions

- Mean squared error: $\ell(x, t; w) = (y - t)^2$

```
torch.nn.MSELoss(size_average=None, reduce=None, reduction='mean')
```

- Minimum log-likelihood: $\ell(x, t; w) = \sum_{(x^{(i)}, t^{(i)}) \in \mathcal{D}} -\log p(t|x)$

- Combined with binary classification: $\ell(x, t; w) = \sum_{(x^{(i)}, t^{(i)}) \in \mathcal{D}} \log(1 + \exp(-t^{(i)} w^T x^{(i)}))$

- Combined with softmax: $\ell(x, t; w) = \sum_{(x^{(i)}, t^{(i)}) \in \mathcal{D}} (-w_{t^{(i)}}^T x + \log \sum_{c \in \{0, \dots, C-1\}} \exp(w_c^T x))$

```
torch.nn.CrossEntropyLoss(weight=None, size_average=None, ignore_index=-100, reduce=None, reduction='mean', label_smoothing=0.0)[source]
```

- Cross Entropy Loss:

- Linear (SVM formulation): $\ell(x, t; w) = \frac{|W[1:]|}{2} + C \sum \max(0, 1 - t^{(i)} \cdot Wx^{(i)})^2$

- Logistic: $\ell(x, t; w) = -t \log y - (1 - t) \log(1 - y)$

Optimizers and torch.optim

In standard gradient descent, the update rule is: $\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha \nabla f(\mathbf{w}_k)$. In *gradient descent with momentum*, we introduce a *velocity term* v_k : $v_{k+1} = \beta v_k - \alpha \nabla f(\mathbf{w}_k)$ and $\mathbf{w}_{k+1} = \mathbf{w}_k + v_{k+1}$ where: α is the learning rate, $\beta \in [0, 1]$ is the momentum coefficient, and v_k is the velocity term.

The following are some useful optimizers provided by the torch.optim library including:

- Stochastic gradient descent

```
torch.optim.SGD(params, lr=0.001, momentum=0, dampening=0, weight_decay=0, nesterov=False, *, maximize=False, foreach=None, differentiable=False, fused=None)[source]
```

PyTorch datasets

Required functions for dataset class:

- `__init__`: The `__init__` method is the constructor for the new dataset.
- `__len__`: The `__len__` method overrides the `len()` function in Python to determine the length of the dataset.
- `__getitem__`: The `__getitem__` method overloads the use of brackets to index items in a dataset.

There are lots of cool dataloader attributes and methods including:

- `batch_size`: number of examples in each batch or call to the dataloader
- `shuffle`: Boolean option to shuffle dataset each pass or *epoch* through the dataset
- `sampler`: *Sampler* object that specifies how data will be extracted from the dataset. For example, the *SubsetRandomSampler* allows us to specify indices within the larger dataset to sample at random.

Other useful equations

- Gradient descent: $\mathbf{w} \leftarrow \mathbf{w} - \alpha \frac{\partial \mathcal{E}}{\partial \mathbf{w}}$
- Closed form solution for linear regression: $W = (X^T X)^{-1} X^T T$
- L2 Regularization with MSE: $L(w) = \|y - Xw\|^2 + \lambda \|w\|_2^2$, closed form linear regression solutions: $W = (X^T X + \lambda I_d)^{-1} X^T y$
- Support Vector Machines - Margins at $WX = 1$ and $WX = -1$, border at $WX = 0$. Margin width = $2/|W|$

Sample Code

Here is a sample, two-dimensional logistic classifier code:

```
import numpy as np
import matplotlib.pyplot as plt
import torch
import torch.nn as nn
from torch.utils.data import Dataset
from torch.utils.data import DataLoader
from torch.utils.data import SubsetRandomSampler

class LogisticRegression(nn.Module):
    def __init__(self, N):
        super().__init__()
        self.w = nn.Parameter(torch.ones(N))
        self.b = nn.Parameter(torch.zeros(1))

    def forward(self, x):
        return 1/(1+torch.exp(-(self.w*x+self.b)))

class TwoClassDataset(Dataset):
    # don't forget the self identifier!
    def __init__(self, N, sigma):
        self.N = N # number of data points per class
        self.sigma = sigma # standard deviation of each class cluster
        self.plus_class = self.sigma*torch.randn(N, 2) + torch.tensor([-1, 1])
        self.negative_class = self.sigma*torch.randn(N, 2) + torch.tensor([1, -1])
        self.data = torch.cat((self.plus_class, self.negative_class), dim=0)
        self.labels = torch.cat((torch.ones(self.N), torch.zeros(self.N)))

    def __len__(self):
        return len(self.labels)

    def __getitem__(self, idx):
        x = self.data[idx]
        y = self.labels[idx]
        return x, y # return input and output pair

N = 100
sigma = 1.5
dataset = TwoClassDataset(N, sigma)
plus_data = dataset.plus_class
negative_data = dataset.negative_class

# create indices for each split of dataset
N_train = 60
N_val = 20
N_test = 20
indices = np.arange(len(dataset))
np.random.shuffle(indices)
train_indices = indices[:N_train]
val_indices = indices[N_train:N_train+N_val]
test_indices = indices[N_train+N_val:]

# create dataloader for each split
batch_size = 8
train_loader = DataLoader(dataset, batch_size=batch_size, sampler=SubsetRandomSampler(train_indices))
val_loader = DataLoader(dataset, batch_size=batch_size, sampler=SubsetRandomSampler(val_indices))
test_loader = DataLoader(dataset, batch_size=batch_size, sampler=SubsetRandomSampler(test_indices))

criterion = nn.BCELoss(reduction='mean') # binary cross-entropy loss, use mean loss
logreg_model = LogisticRegression(2) # initialize model
optimizer = torch.optim.SGD(logreg_model.parameters()) # initialize optimizer

n_epoch = 200 # number of passes through the training dataset
loss_values, train_accuracies, val_accuracies = [], [], []
for n in range(n_epoch):
    epoch_loss, epoch_acc = 0, 0
    for x_batch, y_batch in train_loader:
        optimizer.zero_grad()
        predictions = logreg_model(x_batch.unsqueeze(-1)).squeeze(-1)
        loss = criterion(predictions, y_batch)
        loss.backward()
        optimizer.step()
```