

Sheet 2.3: Non-linear regression (MLP w/ PyTorch modules)

Contents

- Packages & global parameters
- True model & training data
- Defining the MLP using PyTorch's built-in modules
- Preparing the training data
- Training the model
- Outlook: PyTorch layers and utils [optional]

Author: Michael Franke

In this tutorial, we will fit a non-linear regression, implemented as a multi-layer perceptron. We will see how the use of modules from PyTorch's neural network package `torch.nn` helps us implement the model efficiently.

Packages & global parameters

We will need to import the `torch` package for the main functionality. In addition to the previous sheet, In order to have a convenient, we will use PyTorch's `DataLoader` and `Dataset` in order to feed our training data to the model.

```
import torch
import torch.nn as nn
from torch.utils.data import DataLoader
from torch.utils.data import Dataset
import numpy as np
import seaborn as sns
import pandas as pd
import matplotlib.pyplot as plt
import warnings

warnings.filterwarnings("ignore")
```

True model & training data

The "true model" is a constructed non-linear function $y = f(x)$. Here is its definition and a plot to show what the "ground truth" looks like.

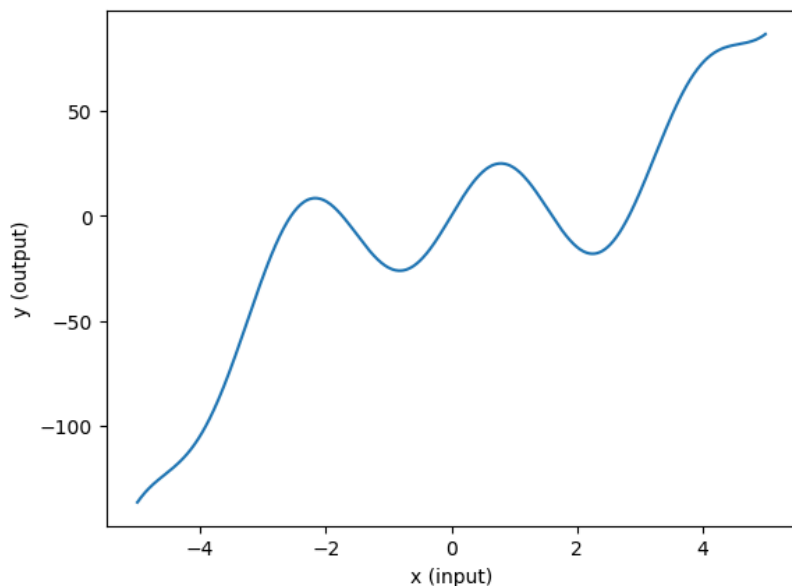
```
#####
## ground-truth model
#####

def goal_fun(x):
    """
    This is a non-linear function that is the "ground truth" for this example.
    """
    return x**3 - x**2 + 25 * np.sin(2 * x)

# create linear sequence (x) and apply goal_fun(y)
x = np.linspace(start=-5, stop=5, num=1000)
y = goal_fun(x)

# plot the function
d = pd.DataFrame({"x (input)": x, "y (output)": y})
sns.lineplot(data=d, x="x (input)", y="y (output)")
plt.show()
```

[Skip to main content](#)



The training data consists of 100 pairs of (x, y) values. Each pair is generated by first sampling an x value from a uniform distribution. For each sampled x , we compute the value of the target function $f(x)$ and add Gaussian noise to it.

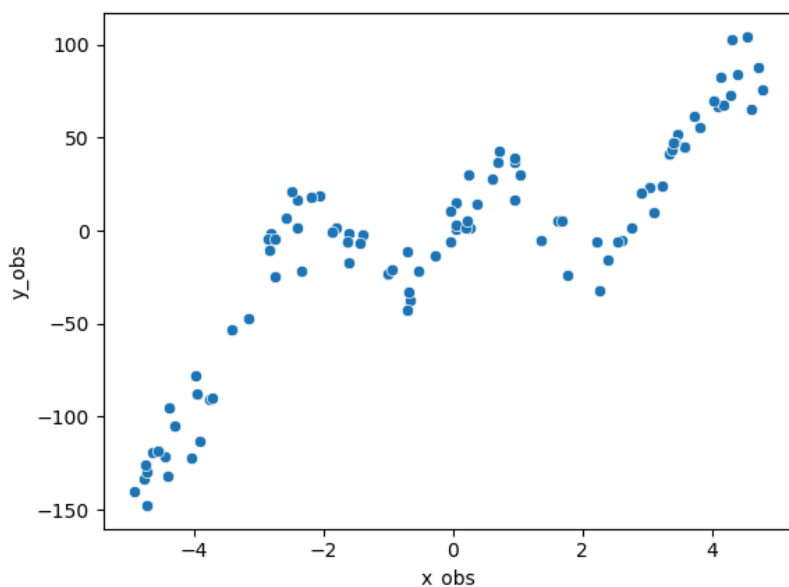
```
#####
## generate training data (with noise)
#####

n_obs = 100 # number of observations

# get noise around y observations
y_normal = torch.distributions.Normal(loc=0.0, scale=10)
y_noise = y_normal.sample([n_obs])

# get observations
x_obs = 10 * torch.rand([n_obs]) - 5 # uniform from [-5,5]
y_obs = goal_fun(x_obs) + y_noise

# plot the data
d = pd.DataFrame({"x_obs": x_obs, "y_obs": y_obs})
sns.scatterplot(data=d, x="x_obs", y="y_obs")
plt.show()
```



Defining the MLP using PyTorch's built-in modules

Our model maps a single scalar x onto another scalar y . We use a 3-layer MLP, each hidden layer with dimension 10:

```
#####
## network dimension parameters
.....
```

[Skip to main content](#)

```
n_input = 1
n_hidden = 10
n_output = 1
```

PyTorch defines a special-purpose class called `nn.Module` from which pre-defined neural networks or custom-made networks inherit the structure and basic functionality. Below, we define our feed-forward neural network as a class extending `nn.Module`. Minimally, we have to define two functions for this to work:

1. the **initialization** function `__init__` which defines which variables (mostly, but not exclusively: parameters) our model has (using `nn.Linear` instantiates a linear layer with all the trainable parameters (weights and biases) implicitly);
2. the **forward pass** which takes the model's input and computes the corresponding prediction given the current parameter values. Instantiating the forward pass is critical because it is designed to also implicitly compute and save gradients of single computations that constitute the neural network with respect to the current input, so that they can be later used for the **backward pass** (i.e., for training).

Since PyTorch allows flexibility in how to define neural network modules, we look at two variants below, one explicit and one more concise. They should implement the exact same model and work the same way eventually.

More explicit definition NN module

Print to PDF ►

more control of the actual forward compared to the concise version below.

```
#####
## set up multi-layer perceptron w/ PyTorch
## -- explicit version --
#####

class MLPexplicit(nn.Module):
    """
    Multi-layer perceptron for non-linear regression.
    """

    def __init__(self, n_input, n_hidden, n_output):
        super(MLPexplicit, self).__init__()
        self.n_input = n_input
        self.n_hidden = n_hidden
        self.n_output = n_output
        self.linear1 = nn.Linear(self.n_input, self.n_hidden)
        self.linear2 = nn.Linear(self.n_hidden, self.n_hidden)
        self.linear3 = nn.Linear(self.n_hidden, self.n_hidden)
        self.linear4 = nn.Linear(self.n_hidden, self.n_output)
        self.relu = nn.ReLU()

    def forward(self, x):
        h1 = self.relu(self.linear1(x))
        h2 = self.relu(self.linear2(h1))
        h3 = self.relu(self.linear3(h2))
        output = self.linear4(h3)
        return output

mlp_explicit = MLPexplicit(n_input, n_hidden, n_output)
```

We can access the current parameter values of this model instance like so:

```
for param in mlp_explicit.parameters():
    print(param.detach().numpy().round(4))
```

have a look at the numbers and make sense of them

weights of the 1st linear layer

```
[[-0.0076]
 [-0.8415]
 [ 0.9536]
 [-0.6464]
 [-0.9964]
 [-0.7288]
 [ 0.1196]
 [-0.462 ]
 [-0.8668]
 [-0.9444]]
[-0.5211  0.0849  0.5   -0.2557  0.9653 -0.4786  0.5713  0.6843  0.2133
 0.342 ]
[[-0.0769 -0.2313  0.0227  0.2852  0.0166  0.2986 -0.0292  0.1179  0.2564
 -0.2944]
 [ 0.1009  0.0552 -0.0447  0.2081  0.2009 -0.033  -0.1969  0.1945 -0.1427
 0.0546]
```

matrix x vector, here is a matrix

a vector of size 10, the bias

[Skip to main content](#)

10 hidden layers

```
0.2865]
[-0.1744 -0.0109 -0.1259 -0.1426 0.0955 -0.3025 -0.1744 -0.2543 0.2068
-0.0964]
[ 0.0934 -0.0013 -0.2462 -0.2037 -0.0994 0.2537 -0.1462 -0.2201 0.1222
0.2095]
[-0.107 -0.253 0.1784 0.1452 -0.0968 0.1959 0.0706 -0.0282 -0.2097
0.1417]
[-0.0523 0.2028 0.2256 0.1229 -0.1682 0.0451 0.0309 0.1407 0.1962
0.0854]
[ 0.2872 -0.0421 0.3001 0.0915 -0.0519 -0.2862 -0.1163 -0.2389 -0.0924
0.2941]
[ 0.2323 0.0592 0.212 -0.071 -0.0909 0.1527 0.0588 -0.2353 -0.0547
0.0211]
[ 0.1491 0.2182 0.2239 0.2984 -0.1006 -0.1116 -0.0044 0.251 -0.0451
0.2208]]
[ 0.039 -0.0528 0.3162 0.0257 0.0663 -0.2846 0.0646 0.2537 -0.1269
0.2371]
[[ 0.3099 0.2536 -0.0446 0.1587 0.0044 -0.255 -0.074 -0.1321 0.2919
-0.1839]
[-0.1293 0.0429 0.2682 0.2615 0.0671 -0.2576 -0.1129 0.0199 -0.0655
-0.043 ]
[-0.1449 0.1796 0.0972 -0.1241 0.2595 -0.151 -0.1908 0.286 -0.2029
0.2647]
[ 0.1574 -0.3025 0.2976 0.1041 0.0947 0.0225 0.0479 -0.0493 -0.2154
-0.2965]
[-0.2934 0.0936 -0.3114 -0.2956 0.2425 0.1305 -0.046 -0.1408 0.0655
0.023 ]
[-0.0155 0.1967 0.2333 -0.179 0. 0.0348 0.074 0.2954 -0.2073
-0.2232]
[ 0.1136 -0.0556 -0.0454 0.315 0.2971 0.0509 0.1394 0.2337 0.0912
-0.256 ]
[-0.1774 0.2738 0.0386 -0.034 0.2706 -0.0666 -0.2664 0.0419 -0.1732
0.0622]
[ 0.2011 0.276 -0.0453 -0.1481 -0.0383 -0.2165 0.1897 -0.1638 0.1007
-0.2027]
[-0.3049 0.2323 0.0017 0.2579 0.0675 0.1464 -0.1149 0.195 -0.1961
-0.2143]]
[ 0.2373 -0.0853 0.016 0.1031 0.1161 -0.1396 0.0505 0.2276 -0.2393
0.1518]
[[ 0.2711 0.1264 -0.2712 -0.2633 0.1861 -0.2758 0.0348 -0.2591 -0.3013
0.2468]]
[0.2455]
```

Exercise 2.3.1: Inspect the model's parameters and their initial values

0. Make sure that you understand what these parameters are. (Hint: We have a 3 layer network, each hidden layer has 10 nodes. We should expect to see slopes (aka weights) and intercepts (biases) somewhere. The order of the presentation in this print-out is the order in which the numbers matter in the computation of the forward pass.)
1. Guess how the weights of the slope matrices are initialized (roughly). Same for the intercept vectors.

More concise definition of NN module

Here is another, more condensed definition of the same NN model, which uses the `nn.Sequential` function to neatly chain components, thus defining the model parameters and the forward pass in one swoop.

```
#####
## set up multi-layer perceptron w/ PyTorch
## -- condensed version --
#####

class MLPCondensed(nn.Module):
    """
    Multi-layer perceptron for non-linear regression.
    """

    def __init__(self, n_input, n_hidden, n_output):
        super().__init__()
        self.layers = nn.Sequential(
            nn.Linear(n_input, n_hidden),
            nn.ReLU(),
            nn.Linear(n_hidden, n_hidden),
            nn.ReLU(),
            nn.Linear(n_hidden, n_hidden),
            nn.ReLU(),
            nn.Linear(n_hidden, n_output),
```

[Skip to main content](#)

```
def forward(self, x):
    return self.layers(x)

mlp_condensed = MLPcondensed(n_input, n_hidden, n_output)
```

Here you can select which of the two models to use in the following.

```
# which model to use from here onwards
# model = mlpExplicit
model = mlp_condensed    the object, the model
```

Preparing the training data

Data pre-processing is a tedious job, but an integral part of machine learning. In order to have a clean interface between data processing and modeling, we would ideally like to have a common data format to feed data into any kind of model. That also makes sharing and reusing data sets much less painful. For this purpose, PyTorch provides two data primitives:

`torch.utils.data.Dataset` and `torch.utils.data.DataLoader`. The class **Dataset** stores the training data (in a reusable format). The class **DataLoader** takes a `dataset` object as input and returns an iterable to enable easy access to the training data.

To define a `Dataset` object, we have to specify two key functions:

1. the `__len__` function, which tells subsequent applications how many data points there are; and [tell system how many data points](#)
2. the `__getitem__` function, which takes an index as input and outputs the data point corresponding to that index. [how to get the next item as thte index](#)

```
#####
## representing training data as a Dataset object
#####

class NonLinearRegressionData(Dataset):
    """
    Custom 'Dataset' object for our regression data.
    Must implement these functions: __init__, __len__, and __getitem__.
    """

    def __init__(self, x_obs, y_obs):
        self.x_obs = torch.reshape(x_obs, (len(x_obs), 1))
        self.y_obs = torch.reshape(y_obs, (len(y_obs), 1))

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

    def __getitem__(self, idx):
        return (x_obs[idx], y_obs[idx])

# instantiate Dataset object for current training data
d = NonLinearRegressionData(x_obs, y_obs)

# instantiate DataLoader
# we use the 4 batches of 25 observations each (full data has 100 observations)
# we also shuffle the data
train_dataloader = DataLoader(d, batch_size=25, shuffle=True)
100 observations in total, do not want to train at the same time as it often does not work.
We only look at a batch of 25 at this moment. And repeat train on different batch, four times
```

We can test the iterable that we create, just to inspect how the data will be delivered later on:

```
for i, data in enumerate(train_dataloader, 0):
    input, target = data
    print(f"In: {input}")
    print(f"Out:{target}\n")
```

```
In: tensor([-4.4495, -2.5819, -4.7650,  0.1632, -2.3451, -1.6135, -1.6223,  0.1976,
          3.7169,  3.3746,  4.5931, -2.8005, -0.0434, -1.8540,  0.0562, -4.4078,
          2.9046,  3.4005, -0.7086, -2.1871,  3.0338, -4.3793, -4.7298,  1.7743,
          -3.4155])
Out: tensor([-121.5500,   7.0789, -133.6013,   2.2198, -21.6659,  -1.4612,
          -5.7734,   1.6199,  61.5975,  43.4716,  65.3842,  -1.3931,
          -6.2188,  -0.4141,   0.6194, -131.9025,  20.4001,  46.9520,
          -43.0037,  18.1146,  23.4007,  -94.9904, -147.4500,  -23.9163,
```

[Skip to main content](#)

```

In: tensor([ 2.2630, -4.0369,  0.3773,  0.2532,  4.5453,  4.2955, -1.3915,  2.7552,
            -0.6925, -2.3938, -1.6076,  1.6834, -2.8546,  3.2283,  3.8123,  3.3403,
            0.9574, -2.0644, -3.9521,  0.0436,  2.2180,  2.3797, -4.7454, -3.7530,
            -2.8264])
Out: tensor([-32.3442, -122.4230,  14.2478,   1.3011, 104.4167, 102.6514,
            -2.1359,   1.3875, -32.5973,  16.5617, -16.8176,   5.5437,
            -4.1630,  23.9920,  55.6051,  41.1025,  38.8970,  18.9164,
            -87.8666, 15.2518,  -6.0591, -15.5863, -125.6087, -90.5321,
            -10.3973])

In: tensor([-3.9640,  0.6017,  0.6880,  0.9468, -1.4264, -1.0024,  4.7179,  3.5653,
            2.5318,  1.0374,  4.7767,  0.7110, -0.6955, -2.7494,  4.0279, -0.5396,
            3.4075,  0.9538, -4.6411, -1.8022, -0.9450,  3.4118,  2.7496,  2.6096,
            -0.6654])
Out: tensor([-77.7501,  27.7172,  36.9245,  16.8191,  -6.6196, -23.5388,
            88.0235,  45.1277,  -6.0365,  30.0066,  76.0195,  42.6423,
            -11.5127, -24.7764,  69.8224, -21.5756,  46.9476,  37.1051,
            -119.4592,  1.4126, -21.1727,  47.0679,  1.6114,  -5.0246,
            -37.5510])

In: tensor([ 2.9079,  4.1315,  0.0575, -2.4921, -2.7523,  4.3927, -3.9023, -2.3946,
            -4.3009, -0.0370,  0.2317,  4.1652,  3.1014, -4.9101, -3.1524,  4.2815,
            4.0928, -0.2689,  1.6177, -3.7121, -4.7293, -4.5396,  3.4614,  0.2180,
            1.3488])
Out: tensor([ 20.1371,  82.1433,   3.2099,  20.7502,  -4.5580,  83.6555,
            -113.4812,  1.6265, -104.7873,  10.2406,  30.3842,  67.8441,
            10.0513, -140.1428, -47.4315,  72.9185,  66.8981, -13.1600,
            5.2063,  -89.8697, -129.4619, -118.5045,  51.6166,   5.1447,
            -5.2186])

```

Training the model

We can now train the model similar to how we did this before. Note that we need to slightly reshape the input data to have the model compute the batched input correctly.

```

#####
## training the model
#####

# Define the loss function and optimizer
loss_function = nn.MSELoss()
optimizer = torch.optim.Adam(model.parameters(), lr=1e-4)
n_train_steps = 50000

# Run the training loop
for epoch in range(0, n_train_steps):

    # Set current loss value
    current_loss = 0.0

    # Iterate over the DataLoader for training data
    for i, data in enumerate(train_dataloader, 0):
        # Get inputs
        inputs, targets = data
        # Zero the gradients
        optimizer.zero_grad()
        # Perform forward pass (make sure to supply the input in the right way)
        outputs = model(torch.reshape(inputs, (len(inputs), 1))).squeeze()
        # Compute loss
        loss = loss_function(outputs, targets)
        # Perform backward pass
        loss.backward()
        # Perform optimization
        optimizer.step()
        # Print statistics
        current_loss += loss.item()

    if (epoch + 1) % 2500 == 0:
        print(f"Loss after epoch {epoch + 1:5d}: {current_loss:.3f}")
        current_loss = 0.0

# Process is complete.
print("Training process has finished.")

y_pred = np.array(
    [model.forward(torch.tensor([o])).detach().numpy() for o in x_obs]
).flatten()

# plot the data
d = pd.DataFrame(
    {"x_obs": x_obs.detach().numpy(), "y_obs": y_obs.detach().numpy(), "y_pred": y_pred}

```

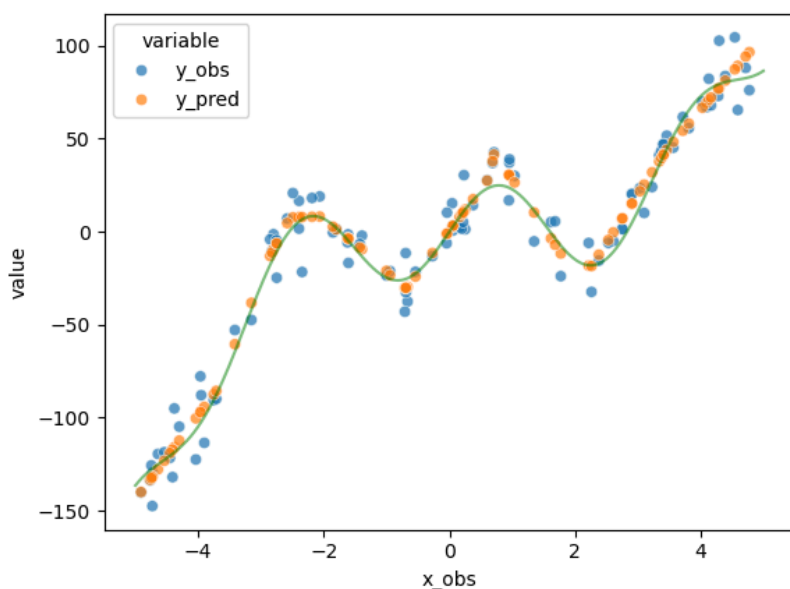
sgd doesn't work in this case.

time consuming, often need to reshape a lot of times
call the function model, see above for reference

make sure the input of the model function is what you expect in the forward

```
sns.scatterplot(data=d_wide, x="x_obs", y="value", hue="variable", alpha=0.7)
x = np.linspace(start=-5, stop=5, num=1000)
y = goal_fun(x)
plt.plot(x, y, color="g", alpha=0.5)
plt.show()
```

```
Loss after epoch 2500: 2956.298
Loss after epoch 5000: 1596.497
Loss after epoch 7500: 740.628
Loss after epoch 10000: 522.447
Loss after epoch 12500: 473.209
Loss after epoch 15000: 450.191
Loss after epoch 17500: 441.309
Loss after epoch 20000: 436.642
Loss after epoch 22500: 433.959
Loss after epoch 25000: 432.140
Loss after epoch 27500: 430.475
Loss after epoch 30000: 428.987
Loss after epoch 32500: 427.638
Loss after epoch 35000: 425.664
Loss after epoch 37500: 424.004
Loss after epoch 40000: 422.595
Loss after epoch 42500: 419.309
Loss after epoch 45000: 417.816
Loss after epoch 47500: 416.479
Loss after epoch 50000: 415.121
Training process has finished.
```



Exercise 2.3.2: Explore the model's behavior

0. [Just for yourself.] Make sure you understand *every line* in this last code block. Ask if anything is unclear.
1. Above we used the DataLoader to train in 4 mini-batches. Change it so that there is only one batch containing all the data. Change the `shuffle` parameter so that data is not shuffled. Run the model and check if you observe any notable differences. Explain what your observations are. (If you do not see anything, explain why you don't. You might pay attention to the results of training.)

Outlook: PyTorch layers and utils [optional]

PyTorch provides vast functionality around neural networks. Next to the `nn.Linear()` layer (a linear layer) and the `nn.ReLU()` (ReLU activation function, implemented as a separate 'layer'), which we have seen in the MLP above, PyTorch implements many other commonly used neural network layers. These different layers can be stacked to create your own custom architectures. Furthermore, PyTorch provides all other building blocks for neural networks like, other activation functions, utilities for 'engineering tricks' that make training neural nets more stable and efficient etc.

The main overview and documentation for all these useful things can be found [here](#). Understanding documentation of useful packages is a core skill for working with code; therefore, the exercise below provides an opportunity to practice this skill and find out more about some treats of PyTorch which are particularly useful for working with language modeling

[Skip to main content](#)

NB: some of the following concepts might be somewhat unclear right now; they should become more clear as the class advances and we learn, e.g., about tokenization. The idea of this and the next part of this sheet is mainly to point your attention to useful things that exist, so that, when needed or when you want to learn more, you will know where to start looking. Therefore, these parts are *optional*.

Exercise 2.3.3: Define functionality of PyTorch utilities

Please look at the documentation references above and, for yourself, write down short (intuitive) definitions of what the following building blocks are (layer or activation function or something else), how / why / what for / where you think they could be used in the context of language modeling.

1. `nn.GELU`:
2. `nn.LogSoftmax`:
3. `nn.GRU`:
4. `nn.GRUCell`:
5. `nn.Transformer`:
6. `nn.Embedding`:
7. `nn.CosineSimilarity`:
8. `nn.CrossEntropyLoss`:
9. `clip_grad_value_`:
10. `weight_norm`:
11. `rnn.pad_sequence`:
12. `rnn.pad_packed_sequence`:

Advanced outlook: PyTorch Autograd and computational graph [optional]

The main work horse behind training neural networks is the chain rule of differentiation. This is the basic rule of calculus that allows exact computation of the gradient of the complicated function that is a neural net with respect to the current input. This gradient, in turn, allows to compute the value by which each weight of the net should be updated in order to improve on the target task (e.g., predicting the next word).

However, computing the gradient even for a simple network like the one above can become difficult. Luckily, PyTorch handles this process behind the scenes with the help of *autograd*, an automatic differentiation technique that computes the gradients with respect to inputs given a *computational graph*. The computational graph is a directed acyclic graph representation of the computations instantiating the neural net, where nodes are operations (e.g., multiplication) and edges represent the data flow between operations. It is created with the definition of the `forward()` method of the model.

If you want to dig deeper into how autograd and PyTorch computational graphs work, here are two great blogposts by PyTorch:

- [Autograd](#)
- [Computational graph](#)

Other blogposts are also worthwhile reading.

Exercise 2.3.4: Understanding PyTorch autograd

1. Try to draw the computational graph (like in Figure 4 of the autograd blogpost) for the function $f(x) = x^3 - x^2$.