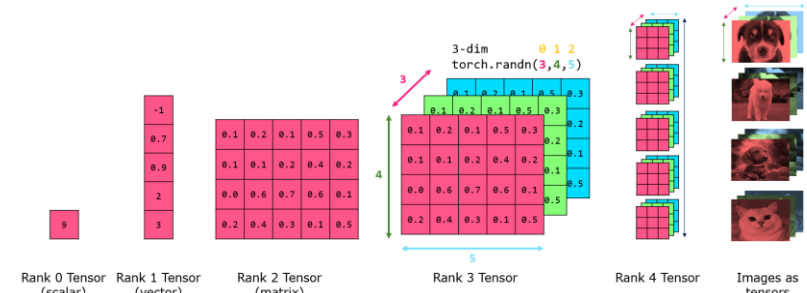


# PyTorch CHEAT SHEET

## Imports

```
import torch # root package
from torch.utils.data import Dataset, DataLoader # dataset representation and loading
```



# N-dim of tensor = number of square brackets [ on the left side  
# [batch\_size, height, width, color\_channels] NHWC # images alternative 1  
# [batch\_size, color\_channels, height, width] NCHW # images alternative 2

## Tensors Creation

```
x = torch.tensor(data, dtype, device) # tensor with no autograd history
x = torch.arange(start, end, step) # sequence
x = torch.randn(*size) # tensor of random numbers from normal distribution N(0,1)
x = torch.ones(*size) # tensor with all 1's [or 0's]
x = torch.tensor(L) # tensor from [nested] list or ndarray L
x = torch.from_numpy(numpy_arr) # tensor from a NumPy array
y = x.clone() # clone of x
with torch.no_grad(): # code wrap that stops autograd from tracking tensor history
    requires_grad=True # arg, when set to True, tracks computation
                        # history for future derivative calculations
```

## Dimensionality

```
x.size() # return tuple-like object of dimensions
x = torch.cat(tensors, dim=0) # concatenates tensors along dim WITHOUT changing the dim of
                               # tensors
y = torch.stack(tensors, dim=0) # stacks a sequence of tensors along a NEW dimension
y = x.view(a, b, ...) # reshapes x into size (a, b, ...)
y = x.view(-1, a) # reshapes x into size (b, a) for some b
y = x.transpose(a, b) # swaps dimensions a and b
y = x.permute(*dims) # Returns a view of the original input with its dimensions
                     # permuted (rearranged) to dims
y = x.unsqueeze(dim) # tensor with added axis
y = x.unsqueeze(dim=2) # (a, b, c) tensor -> (a, b, 1, c) tensor
y = x.squeeze() # removes all dimensions of size 1 (a, 1, b, 1) -> (a, b)
y = x.squeeze(dim=1) # removes specified dimension of size 1 (a, 1, b, 1) -> (a, b, 1)
y = x.reshape(shape) # Reshapes input to shape (if compatible)
```

## Math

```
# Algebra
ret = A * B # element-wise multiplication
ret = A.mm(B) # matrix multiplication / dot product
ret = A.mv(x) # matrix-vector multiplication
x = x.t() # matrix transpose
```

```
torch.abs(tensor)
torch.add(tensor, tensor2) # or tensor*scalar
torch.div(tensor, tensor2) # or tensor/scalar
torch.mult(tensor, tensor2) # or tensor*scalar
torch.sub(tensor, tensor2) # or tensor-scalar
torch.ceil(tensor)
torch.floor(tensor)
torch.remainder(tensor, divisor) #or torch.fmod()
torch.sqrt(tensor)
```

## Torchscript and JIT

```
torch.jit.trace() # takes your module or function and an example
                  # data input, and traces the computational steps
                  # that the data encounters as it progresses through the model

@script # decorator used to indicate data-dependent
        # control flow within the code being traced
```

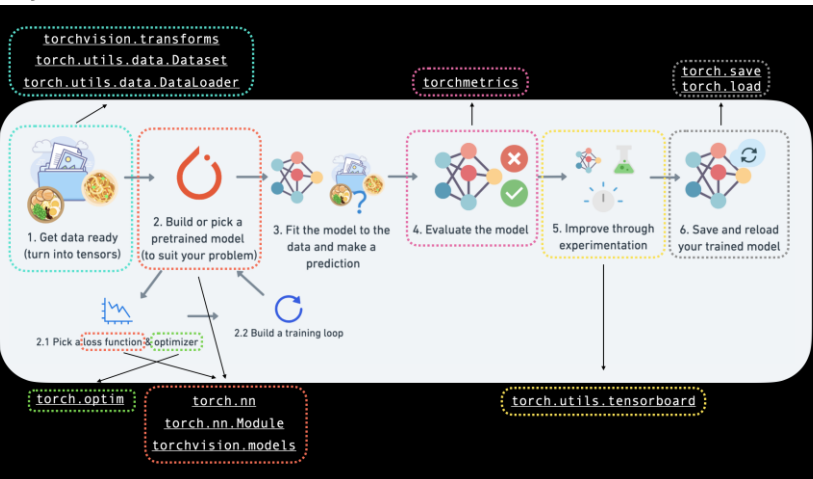
## ONNX

```
torch.onnx.export(model, dummy_data, xxxx.proto) # exports an ONNX formatted
                                                  # model using a trained model, dummy
                                                  # data and the desired file name

model = onnx.load("alexnet.proto") # load an ONNX model
onnx.checker.check_model(model) # check that the model
                                # IR is well formed

onnx.helper.printable_graph(model.graph) # print a human readable
                                         # representation of the graph
```

## PyTorch Workflow



## Data Utilities

```
Dataset # abstract class representing dataset
TensorDataset # labelled dataset in the form of tensors
Concat Dataset # concatenation of Datasets

DataLoader(dataset, batch_size=1, ...) # loads data batches agnostic
                                         # of structure of individual data points

sampler.Sampler(dataset, ...) # abstract class dealing with
                               # ways to sample from dataset

sampler.XSampler where ... # Sequential, Random, SubsetRandom,
                           # WeightedRandom, Batch, Distributed
```

## Vision

```
# Base computer vision library
import torchvision

# Other components of TorchVision (premade datasets, pretrained models and image transforms)
from torchvision import datasets, models, transforms
```

## Text

```
# Base text and natural language processing library
import torchtext

# Other components of TorchText (premade datasets, pretrained models and text transforms)
from torchtext import datasets, models, transforms
```

## Audio and Speech

```
# Base audio and speech processing library
import torchaudio

# Other components of TorchAudio (premade datasets, pretrained models and text transforms)
from torchaudio import datasets, models, transforms
```

## Recommendation systems

```
# Base recommendation system library
import torchrec

# Other components of TorchRec
from torchrec import datasets, models
```

## GPU Usage

```
torch.cuda.is_available() # check for cuda
x = x.cuda() # move x's data from CPU to GPU
              # and return new object

x = x.cpu() # move x's data from GPU to CPU
            # and return new object

if torch.cuda.is_available(): # Setup device-agnostic code
    device = "cuda" # NVIDIA GPU
elif torch.backends.mps.is_available():
    device = "mps" # Apple GPU
else:
    device = "cpu"

net.to(device) # recursively convert their parameters
               # and buffers to device specific tensors

x = x.to(device) # copy your tensors to a device (gpu, cpu)
```

## Neural Network API

```
import torch.autograd as autograd # computation graph
from torch import Tensor # tensor node in the computation graph
import torch.nn as nn # neural networks
import torch.nn.functional as F # layers, activations and more
import torch.optim as optim # optimizers e.g. gradient descent, ADAM, etc.
from torch.jit import script, trace # hybrid frontend decorator and tracing jit
```

## Deep Learning

```
nn.Linear(m, n) # fully connected layer from
                # m to n units

nn.ConvXd(m, n, s) # X dimensional conv layer from
                  # m to n channels where X∈{1,2,3}
                  # and the kernel size is s

nn.MaxPoolXd(s) # X dimension pooling layer
                # (notation as above)

nn.BatchNormXd # batch norm layer
nn.RNN/LSTM/GRU # recurrent layers
nn.Dropout(p=0.5, inplace=False) # dropout layer for any dimensional input
nn.Dropout2d(p=0.5, inplace=False) # 2-dimensional channel-wise dropout
nn.Embedding(num_embeddings, embedding_dim) # (tensor-wise) mapping from
                                           # indices to embedding vectors
```

## Model Building

```
class NeuralNet(nn.Module):
    def __init__(self, input_size, hidden_size, num_classes):
        super().__init__()
        self.fc1 = nn.Linear(input_size, hidden_size)
        self.relu = nn.ReLU()
        self.fc2 = nn.Linear(hidden_size, num_classes)

    def forward(self, x):
        out = self.fc1(x)
        out = self.relu(out)
        out = self.fc2(out)
        return out

model = NeuralNet(input_size, hidden_size, num_classes).to(device)
```

## Loss Functions

```
nn.X # where X is L1Loss, MSELoss, CrossEntropyLoss
      # CTCross, NLLoss, PoissonNLLoss,
      # KLDivLoss, BCELoss, BCEWithLogitsLoss,
      # MarginRankingLoss, HingeEmbeddingLoss,
      # MultiLabelMarginLoss, SmoothL1Loss,
      # SoftMarginLoss, MultiLabelSoftMarginLoss,
      # CosineEmbeddingLoss, MultiMarginLoss,
      # or TripletMarginLoss
```

## Activation Functions

```
nn.X # where X is ReLU, ReLU6, ELU, SELU, PReLU, LeakyReLU,
      # RReLU, CELU, GELU, Threshold, Hardshrink, HardTanh,
      # Sigmoid, LogSigmoid, Softplus, SoftShrink,
      # Softsign, Tanh, TanhShrink, Softmin, Softmax,
      # Softmax2d, LogSoftmax or AdaptiveSoftmaxWithLoss
```

## Optimizers

```
opt = optim.X(model.parameters(), ...) # create optimizer
opt.step() # update weights
opt.X # where X is SGD, Adadelta, Adagrad, Adam,
       # AdamW, SparseAdam, Adamax, ASGD,
       # LBFGS, RMSprop or Rprop
```

## Learning rate scheduling

```
scheduler = optim.X(optimizer, ...) # create lr scheduler
scheduler.step() # update lr after optimizer updates weights
optim.lr_scheduler.X # where X is LambdaLR, MultiplicativeLR,
                    # StepLR, MultiStepLR, ExponentialLR,
                    # CosineAnnealingLR, ReduceLROnPlateau, CyclicalLR,
                    # OneCycleLR, CosineAnnealingWarmRestarts,
```

## Distributed Training

```
import torch.distributed as dist # distributed communication
from torch.multiprocessing import Process # memory sharing processes
```

Ahmad Alismail, [in](https://www.linkedin.com/in/ahmadalismail1) [linkedin.com/in/ahmadalismail1](https://www.linkedin.com/in/ahmadalismail1)

Last updated: 11/2022

PyTorch Workflow: Daniel Bourke <https://github.com/mrdbourke/pytorch-deep-learning>

## Get data ready

- What shapes are my inputs and what shapes are my outputs? (All samples vs. Single sample)
- PyTorch prefers to work with PyTorch tensors -> Turn data into tensors
- Split our data into training and test sets

## Load our data into PyTorch

- Turn int into tensors (e.g., `torchvision.transforms`)
- turn our raw images into Dataset's (features mapped to labels or X's mapped to y's) through `torch.utils.data.Dataset`
- Turn our Dataset subsequently into `DataLoader` through `torch.utils.data.DataLoader`

## DataLoader

- `DataLoader` combines dataset and sampler
- It helps load data into a model. For training and for inference.
- It turns a large **Dataset** into a Python iterable of smaller chunks.
- These smaller chunks are called batches or mini-batches and can be set by the `batch_size` parameter.
- Why do this?
- Because it's more computationally efficient.

## Going Modular

- One way to go from notebooks to scripts
- Train a model built in a jupyter notebook with one line of code on the command line `!python train.py`
- Create folder to store scripts using `os.makedirs()`
- Create the script using `%writefile folder_name/file_name.py`

## Build or pick a pretrained Model

- Setting up device agnostic code (so our model can run on CPU or GPU if it's available)
- Costructing a model by subclassing `nn.Module`
- Create layers capable of handling X and y (considering input and output shapes)
- Define a forward method (make use of operator fusion to improve performance on GPU):  
`return self.layer_3(self.relu(self.layer_2((x)))`
- Create an instance of the model and send it to target device
- Create a loss function (a.k.a loss criterion)
- Create an optimizer

## PyTorch Training Loop

- Put data on the available device `X_train.to(device)` (without this error will happen!)
- Loop through epochs
  - Loop through training batches, perform **training steps**, caculate loss per batch.
- Put model in training mode
- **Forward pass** - The model goes through all of the training data once, performing its `forward()` function calculations (`model(x_train)`)
- **Calculate the loss** - The model's output/predictions (logits for classification if the loss function has a layer to convert logits to probabilities) are compared to the ground truth (labels) and evaluated to see how wrong they are (`loss = loss_fn(y_pred, y_train)`).
- **Zero gradients** - The optimizers gradients are set to zero (they are accumulated by default) so they can be recalculated for the specific training step (`optimizer.zero_grad()`).
- **Perform backpropagation on the loss** - Computes the gradient of the loss with respect for every model parameter to be updated (each parameter with `requires_grad=True`). This is known as backpropagation, hence "backwards" (`loss.backward()`).
- **Step the optimizer (gradient descent)** - Update the parameters with `requires_grad=True` with respect to the loss gradients in order to improve them (`optimizer.step()`).

## PyTorch Testing Loop

- Put the model in evaluation mode
- Turn on `torch.inference_mode()` context manager to disable functionality such as gradient tracking for inference (since gradient tracking is not needed for inference)
- All predictions should be made with objects on the same device (e.g. data and model on **GPU** only or data and model on **CPU** only), i.e., `model.to(device)` and `X_test = X_test.to(device)`
- Pass the test data through the model. The model will generate logits (the raw outputs). Convert **logits** -> prediction **probabilities** (with `sigmoid/softmax`) -> predictions **labels** (with `argmax(dim=1)`)
- NOTE: call `y_pred.cpu()` on your target tensor to return a copy of your target tensor on the CPU because some libraries aren't capable of using data that is stored on GPU

## Top PyTorch and DL Errors

- **Wrong datatypes:** Your Model expected `torch.float32` when your data is `torch.uint8`
- **Wrong data shapes:** Your model expected `[batch_size, color_channels, height, width]` when your data is `[color_channels, height, width]`
- **Wrong devices:** Your model is on the GPU but your data is on the CPU.

## Prevent Overfitting

- [Interpreting Learning Curve](#)
- Get more data
- Use regularization
- Simplify your model
- Use data augmentation
- Use transfer learning
- Use dropout layers
- Use learning rate decay
- Use early stopping

## Prevent Underfitting

- Add more layers/units to your model
- Tweak the learning rate
- Use transfer learning
- Train for longer
- Use less regularization

```
1 # Create a linear regression model in PyTorch
2 class LinearRegressionModel(nn.Module):
3     def __init__(self):
4         super().__init__()
5
6         # Initialize model parameters
7         self.weights = nn.Parameter(torch.randn(1,
8         requires_grad=True,
9         dtype=torch.float
10        ))
11
12        self.bias = nn.Parameter(torch.randn(1,
13        requires_grad=True,
14        dtype=torch.float
15        ))
16
17    # forward() defines the computation in the model
18    def forward(self, x: torch.Tensor) -> torch.Tensor:
19        return self.weights * x + self.bias
20
```

Subclass `nn.Module` (this contains all the building blocks for neural networks)

Initialise **model parameters** to be used in various computations (these could be different layers from `torch.nn`, single parameters, hard-coded values or functions)

`requires_grad=True` means PyTorch will track the gradients of this specific parameter for use with `torch.autograd` and gradient descent (for many `torch.nn` modules, `requires_grad=True` is set by default)

Any subclass of `nn.Module` needs to override `forward()` (this defines the forward computation of the model)

# PyTorch training loop

```
1 # Setup empty lists to keep track of model progress:
2 epoch_count = []
3 train_loss_values = []
4 test_loss_values = []
5
6 # Pass the data through the model for a number of epochs (e.g. 100) epochs:
7 for epoch in range(epochs):
8
9     ### Training loop code here ###
10
11     ### Testing starts ###
12
13     # Put the model in evaluation mode
14     model.eval()
15
16     # Turn on inference mode context manager
17     with torch.inference_mode():
18         # 1. Forward pass on test data
19         test_pred = model(X_test)
20
21         # 2. Calculate loss on test data
22         test_loss = loss_fn(test_pred, y_test)
23
24     # Print out what's happening every 10 epochs
25     if epoch % 10 == 0:
26         epoch_count.append(epoch)
27         train_loss_values.append(loss)
28         test_loss_values.append(test_loss)
29         print(f"Epoch: {epoch} | MAE Train Loss: {loss} | MAE Test Loss: {test_loss}")

```

Create empty lists for storing useful values (helpful for tracking model progress)

Tell the model we want to **evaluate** rather than train (this turns off functionality used for training but not evaluation)

Turn on `torch.inference_mode()` context manager to disable functionality such as gradient tracking for inference (gradient tracking not needed for inference)

Pass the test data through the model (this will call the model's implemented `forward()` method)

Calculate the **test loss value** (how wrong the model's predictions are on the test dataset, lower is better)

Display **information outputs** for how the model is doing during training/testing every ~10 epochs (note: what gets printed out here can be adjusted for specific problems)

Note: all of this can be turned into a function

```
1 # Pass the data through the model for a number of epochs (e.g. 100)
2 for epoch in range(epochs):
3
4     # Put model in training mode (this is the default state of a model)
5     model.train()
6
7     # 1. Forward pass on train data using the forward() method inside
8     y_pred = model(X_train)
9
10    # 2. Calculate the loss (how different are the model's predictions to the true values)
11    loss = loss_fn(y_pred, y_true)
12
13    # 3. Zero the gradients of the optimizer (they accumulate by default)
14    optimizer.zero_grad()
15
16    # 4. Perform backpropagation on the loss
17    loss.backward()
18
19    # 5. Progress/step the optimizer (gradient descent)
20    optimizer.step()

```

Pass the data through the model for a number of **epochs** (e.g. 100 for 100 passes of the data)

Pass the data through the model, this will perform the **forward()** method located within the model object

Calculate the **loss value** (how wrong the model's predictions are)

Zero the **optimizer gradients** (they accumulate every epoch, zero them to start fresh each forward pass)

Perform **backpropagation** on the loss function (compute the gradient of every parameter with `requires_grad=True`)

Step the **optimizer** to update the model's parameters with respect to the gradients calculated by `loss.backward()`

Note: all of this can be turned into a function











## General

PyTorch is a open source machine learning framework. It uses **torch.Tensor** – multi-dimensional matrices – to process. A core feature of neural networks in PyTorch is the autograd package, which provides automatic derivative calculations for all operations on tensors.

import torch	Root package	torch.randn(*size)	Create random tensor
import torch.nn as nn	Neural networks	torch.Tensor(L)	Create tensor from list
from torchvision import datasets, models, transforms	Popular image datasets, architectures & transforms	tnsr.view(a,b, ...)	Reshape tensor to size (a, b, ...)
import torch.nn.functional as F	Collection of layers, activations & more	requires_grad=True	tracks computation history for derivative calculations

## Layers

	nn.Linear(m, n): Fully Connected layer (or dense layer) from m to n neurons		nn.ConvXd(m, n, s): X-dimensional convolutional layer from m to n channels with kernel size s; $X \in \{1, 2, 3\}$
	nn.Flatten(): Flattens a contiguous range of dimensions into a tensor		nn.MaxPoolXd(s): X-dimensional pooling layer with kernel size s; $X \in \{1, 2, 3\}$
	nn.Dropout(p=0.5): Randomly sets input elements to zero during training to prevent overfitting		nn.BatchNormXd(n): Normalizes a X-dimensional input batch with n features; $X \in \{1, 2, 3\}$
	nn.Embedding(m, n): Lookup table to map dictionary of size m to embedding vector of size n		nn.RNN/LSTM/GRU: Recurrent networks connect neurons of one layer with neurons of the same or a previous layer

torch.nn offers a bunch of other building blocks.

A list of state-of-the-art architectures can be found at <https://paperswithcode.com/sota>.

## Load data

A dataset is represented by a class that inherits from **Dataset** (resembles a list of tuples of the form (features, label)).

**DataLoader** allows to load a dataset without caring about its structure.




Usually the dataset is split into training (e.g. 80%) and test data (e.g. 20%).

```
1 from torch.utils.data
2 import Dataset, TensorDataset,
3   DataLoader, random_split
4
5 train_data, test_data =
6   random_split(
7     TensorDataset(inps, tgts),
8     [train_size, test_size]
9   )
10
11 train_loader =
12   DataLoader(
13     dataset=train_data,
14     batch_size=16,
15     shuffle=True)
```

## Activation functions

Common activation functions include **ReLU**, **Sigmoid** and **Tanh**, but there are other activation functions as well.

nn.ReLU() creates a nn.Module for example to be used in Sequential models. F.relu() is just a call of the ReLU function e.g. to be used in the forward method.

	nn.ReLU() or F.relu() Output between 0 and $\infty$ , most frequently used activation function
	nn.Sigmoid() or F.sigmoid() Output between 0 and 1, often used for predicting probabilities
	nn.Tanh() or F.tanh() Output between -1 and 1, often used for classification with two classes

## Define model

There are several ways to define a neural network in PyTorch, e.g. with **nn.Sequential** (a), as a class (b) or using a combination of both.

```
model = nn.Sequential(
  nn.Conv2D(1, 1, 1)
  nn.MaxPool2D(1)
  nn.ReLU()
  nn.Flatten()
  nn.Linear(1, 1)
)
```

```
class Net(nn.Module):
  def __init__():
    super(Net, self).__init__()

    self.conv = nn.Conv2D(1, 1, 1)
    self.pool = nn.MaxPool2D(1)
    self.fc = nn.Linear(1, 1)

  def forward(self, x):
    x = self.pool(F.relu(self.conv(x)))
    x = x.view(-1, 1)
    x = self.fc(x)
    return x
model = Net()
```

## Train model

### LOSS FUNCTIONS

PyTorch already offers a bunch of different loss functions, e.g.:

nn.L1Loss	Mean absolute error
nn.MSELoss	Mean squared error (L2Loss)
nn.CrossEntropyLoss	Cross entropy, e.g. for single-label classification or unbalanced training set
nn.BCELoss	Binary cross entropy, e.g. for multi-label classification or autoencoders

### OPTIMIZATION (torch.optim)

Optimization algorithms are used to update weights and dynamically adapt the learning rate with gradient descent, e.g.:

optim.SGD	Stochastic gradient descent
optim.Adam	Adaptive moment estimation
optim.Adagrad	Adaptive gradient
optim.RMSProp	Root mean square prop

```
1 correct = 0 # correctly classified
2 total = 0 # classified in total
3
4 model.eval()
5 with torch.no_grad():
6   for data in test_loader:
7     inputs, labels = data
8     outputs = model(inputs)
9     _, predicted = torch.max(outputs.data, 1)
10    total += labels.size(0) # batch size
11    correct += (predicted==labels)
12               .sum().item()
13
14 print('Accuracy: %s' % (correct/total))
```

## Save/Load model

model = torch.load('PATH')	Load model
torch.save(model, 'PATH')	Save model

It is common practice to save only the model parameters, not the whole model using model.state\_dict()

```
1 torch.save(model.state_dict(), 'params.ckpt')
2 model.load_state_dict(
3   torch.load('params.ckpt'))
```

## GPU Training

device = torch.device('cuda:0' if torch.cuda.is\_available() else 'cpu')

If a GPU with CUDA support is available, computations are sent to the GPU with ID 0 using model.to(device) or inputs, labels = data[0].to(device), data[1].to(device).

```
1 import torch.optim as optim
2
3 # Define loss function
4 loss_fn = nn.CrossEntropyLoss()
5
6 # Choose optimization method
7 optimizer = optim.SGD(model.parameters(),
8   lr=0.001, momentum=0.9)
9
10 # Loop over dataset multiple times (epochs)
11 for epoch in range(2):
12   model.train() # activate training mode
13   for i, data in enumerate(train_loader, 0):
14     # data is a batch of [inputs, labels]
15     inputs, labels = data
16
17     # zero gradients
18     optimizer.zero_grad()
19
20     # calculate outputs
21     outputs = model(inputs)
22     # calculate loss & backpropagate error
23     loss = loss_fn(outputs, labels)
24     loss.backward()
25     # update weights & learning rate
26     optimizer.step()
```

## Evaluate model

The evaluation examines whether the model provides satisfactory results on previously withheld data. Depending on the objective, different metrics are used, such as accuracy, precision, recall, F1, or BLEU.

model.eval()	Activates evaluation mode, some layers behave differently
torch.no_grad()	Prevents tracking history, reduces memory usage, speeds up calculations