

★ Get unlimited access to all of Medium. [Become a member](#)



A batch too large: Finding the batch size that fits on GPUs

A simple function to identify the batch size for your PyTorch model that can fill the GPU memory



Bryan M. Li · [Follow](#)

Published in Towards Data Science

5 min read · Oct 19, 2022



Listen



Share

... More

I am sure many of you had the following painful experience: you start multiple ML experiments on your GPUs to train overnight and when you came back to check 10 hours later, you realize that the progress bar has barely moved due to hardware underutilization, or worse, all the experiments failed due to out-of-memory (OOM) error. In this mini-guide, we will implement an automated method to find the batch size for your PyTorch model that can utilize the GPU memory sufficiently without causing OOM!



Photo by [Nana Dua](#) on [Unsplash](#)

On top of the model architecture and the number of parameters, the batch size is the most effective hyperparameter to control the amount of GPU memory an experiment uses. The proper method to find the optimal batch size that can fully utilize the accelerator is via GPU profiling, a process to monitor processes on the computing device. Both [TensorFlow](#) and [PyTorch](#) provide detailed guides and tutorials on how to perform profiling in their framework. In addition, the batch size can greatly affect the performance of the model. For instance, a large batch size can lead to poor generalization, check out this blog post by [Kevin Shen](#) on the [effect of batch size on training dynamics](#) if you are interested in this topic.

Nevertheless, if you simply want to train a model to test an idea, profiling or performing a hyperparameter search to find the best batch size might be overkill, especially in the early stage of the project. A common approach to find the value that allows you to fit your model without OOM is to train the model with a small batch size while monitoring the GPU utilization using tools like `nvidia-smi` or [nvidia-top](#). You then increase the value if the model is underutilizing the GPU memory, and repeat the

process until you hit the memory capacity. However, this manual process can be time-consuming. More annoyingly, when you have to run experiments on different GPUs with varying memory sizes then you have to repeat the same process for each device. Luckily, we can convert this tedious iterative process into code and run it before the actual experiment so that you know your model won't cause an OOM.

The idea is very simple:

1. Initialize your model.
2. Set batch size to 2 (for BatchNorm)
3. Create dummy data that has the sample shape as the real data.
4. Train the model for n steps (both forward and backward passes).
5. If the model ran without an error, then increase the batch size and go to Step 3. If OOM is raised (i.e. `RuntimeError` in PyTorch) then set the batch size to the previous value and terminate.
6. Return the final batch size.

To put this into code

```
1  def get_batch_size(  
2      model: nn.Module,  
3      device: torch.device,  
4      input_shape: t.Tuple[int, int, int],  
5      output_shape: t.Tuple[int],  
6      dataset_size: int,  
7      max_batch_size: int = None,  
8      num_iterations: int = 5,  
9  ) -> int:  
10     model.to(device)  
11     model.train(True)  
12     optimizer = torch.optim.Adam(model.parameters())  
13  
14     batch_size = 2  
15     while True:  
16         if max_batch_size is not None and batch_size >= max_batch_size:  
17             batch_size = max_batch_size  
18             break  
19         if batch_size >= dataset_size:  
20             batch_size = batch_size // 2  
21             break  
22         try:  
23             for _ in range(num_iterations):  
24                 # dummy inputs and targets  
25                 inputs = torch.rand(*(batch_size, *input_shape), device=device)  
26                 targets = torch.rand(*(batch_size, *output_shape), device=device)  
27                 outputs = model(inputs)  
28                 loss = F.mse_loss(targets, outputs)  
29                 loss.backward()  
30                 optimizer.step()  
31                 optimizer.zero_grad()  
32                 batch_size *= 2  
33             except RuntimeError:  
34                 batch_size //= 2  
35                 break  
36     del model, optimizer  
37     torch.cuda.empty_cache()  
38     return batch_size
```

find_batch_size.py hosted with ♥ by GitHub

[view raw](#)

As you can see, this function has 7 arguments:

- `model` — the model you want to fit, note that the model will be deleted from memory at the end of the function.
- `device` — `torch.device` which should be a CUDA device.
- `input_shape` — the input shape of the data.
- `output_shape` — the expected output shape of the model.
- `dataset_size` — the size of your dataset (we wouldn't want to continue the search when the batch size is already larger than the size of the dataset).
- `max_batch_size` — an optional argument to set the maximum batch size to use.
- `num_iterations` — the number of iterations to update the model before increasing the batch size, default to 5.

Let's quickly go through what's happening in the function. We first load the model to the GPU, initialize Adam optimizer, and set the initial batch size to 2 (you can start with a batch size of 1 if you are not using BatchNorm). We can then begin the iterative process. First, we check if the current batch size is larger than the size of the dataset or the maximum desired batch size, if so, we break the loop. Otherwise, we create dummy inputs and targets, move them to GPU and fit the model. We train the model for 5 steps to ensure neither forward nor backward pass causes OOM. If everything is fine, we multiply the batch size by 2 and re-fit the model. If OOM occurs during the above steps, then we reduce the batch size by a factor of 2 and exit the loop. Finally, we clear the model and optimizer from memory and return the final batch size. That's it!

Note that, instead of simply dividing the batch size by 2 if the case of OOM, one could continue to search for the optimal value (i.e. binary search the batch size, set batch size to the mid-point between the breaking and last working value, and continue to Step 3.) to find the batch size that fit perfectly to the GPU. However, keep in mind that PyTorch/TensorFlow or other processes might request more GPU memory in the middle of an experiment and you risk OOM, I hence prefer having some wiggle room.

Now let's put this function into use. Here we fit the ResNet50 on 1,000 train synthetic images of size (3, 224, 224) generated by FakeData Datasets. Briefly, we first call

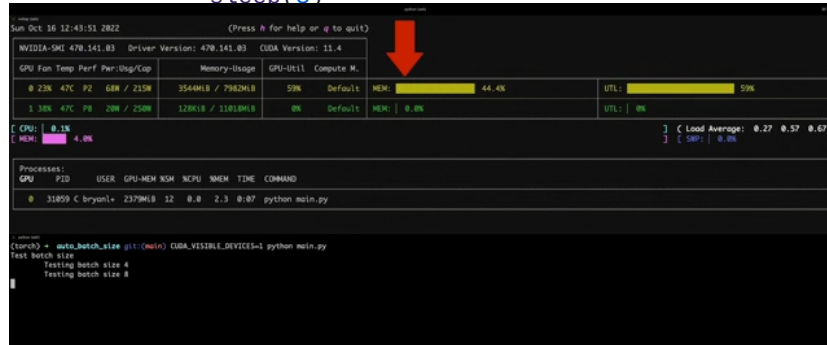
`get_batch_size=(model=ResNet(), input_shape=IMAGE_SHAPE, output_shape=(NUM_CLASSES,), dataset_size=DATASET_SIZE)` to get the batch size that can fill the GPU memory sufficiently. Then we can initialize the model and DataLoaders, and train the model like you normally do!

```
1  import torch
2  import typing as t
3  import torch.nn as nn
4  from tqdm import tqdm
5  import torch.optim as optim
6  import torch.nn.functional as F
7  from torch.utils.data import DataLoader
8  from torchvision import datasets, transforms
9  from torchvision.models import resnet50, ResNet50_Weights
10
11
12  from time import sleep
13
14  # dataset information
15  IMAGE_SHAPE = (3, 224, 224)
16  NUM_CLASSES = 100
17  DATASET_SIZE = 1000
18
19
20  def get_batch_size(
21      model: nn.Module,
22      device: torch.device,
23      input_shape: t.Tuple[int, int, int],
24      output_shape: t.Tuple[int],
25      dataset_size: int,
26      max_batch_size: int = None,
27      num_iterations: int = 5,
28  ) -> int:
29      model.to(device)
30      model.train(True)
31      optimizer = torch.optim.Adam(model.parameters())
32
33      print("Test batch size")
34      batch_size = 2
35      while True:
36          if max_batch_size is not None and batch_size >= max_batch_size:
37              batch_size = max_batch_size
38              break
39          if batch_size >= dataset_size:
40              batch_size = batch_size // 2
41              break
42          try:
43              for _ in range(num_iterations):
44                  # dummy inputs and targets
45                  inputs = torch.rand(*(batch_size * input_shape), device=device)
```

```

45     inputs = torch.rand( (batch_size, *input_shape), device=device)
46     targets = torch.rand(* (batch_size, *output_shape), device=device)
47     outputs = model(inputs)
48     loss = F.mse_loss(targets, outputs)
49     loss.backward()
50     optimizer.step()
51     optimizer.zero_grad()
52     batch_size *= 2
53     print(f"\tTesting batch size {batch_size}")
54     sleep(3)

```



8.80 s

SD

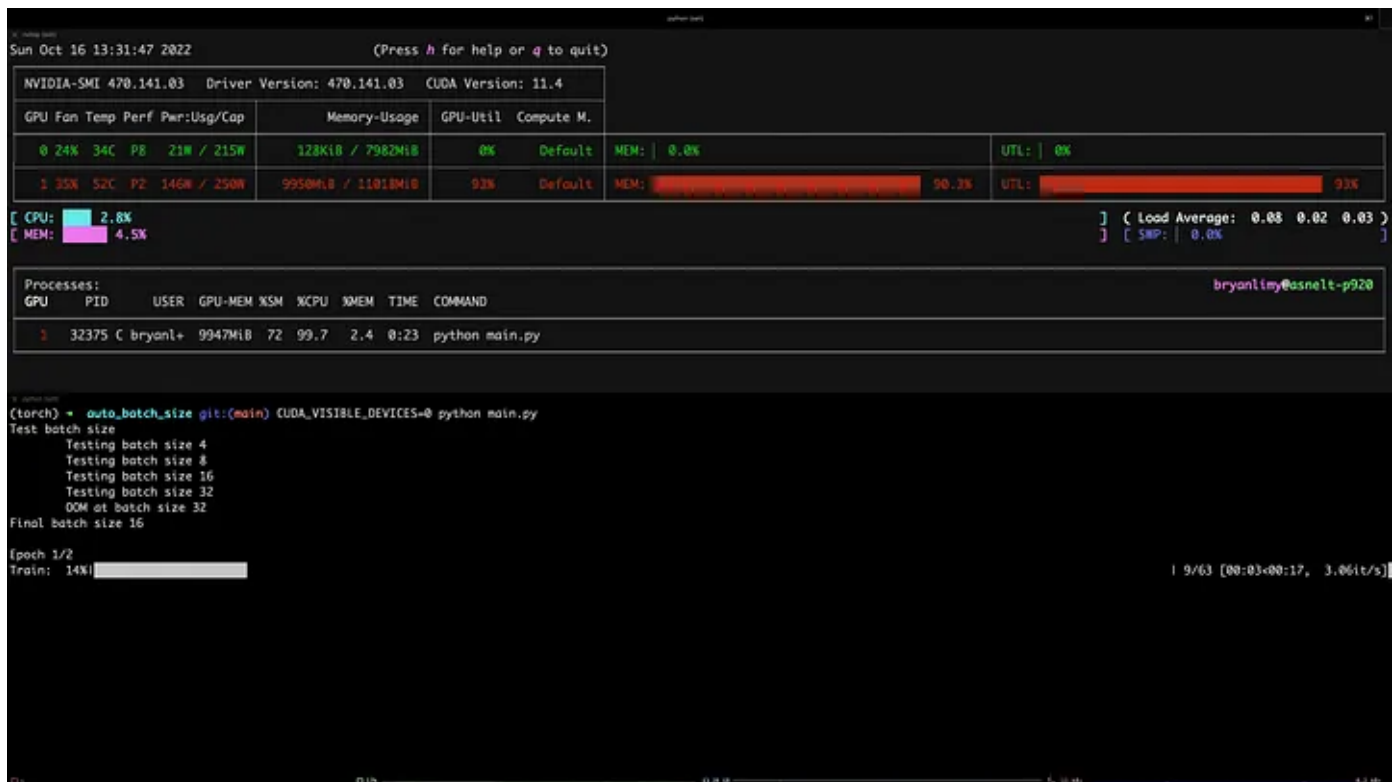
7.9K views



```

67     datasets.FakeData(
68         size=DATASET_SIZE,
69         image_size=IMAGE_SHAPE,
70         num_classes=NUM_CLASSES,
71         transform=transforms.Compose([transforms.ToTensor()]),
72     ),
73     batch_size=batch_size,
74     shuffle=True,
75     num_workers=num_workers,
76 )
77     test_ds = DataLoader(
78         datasets.FakeData(
79             size=200,
80             image_size=IMAGE_SHAPE,
81             num_classes=NUM_CLASSES,
82             transform=transforms.Compose([transforms.ToTensor()]),
83         ),
84         batch_size=batch_size,
85         num_workers=num_workers,
86     )
87     return train_ds, test_ds
88
89

```

```

108     optimizer: torch.optim,
109     train_ds: DataLoader,
110     device: torch.device,
111 ):
112     model.train()
113     train_loss, correct = 0, 0
114     for batch_idx, (data, target) in enumerate(tqdm(train_ds, desc="Train")):
115         data, target = data.to(device), target.to(device)
116         optimizer.zero_grad()
117         output = model(data)
118         loss = F.nll_loss(output, target)
119         train_loss += loss.item()
120         loss.backward()
121         optimizer.step()
122         pred = output.max(1, keepdim=True)[1]
123         correct += pred.eq(target.view_as(pred)).sum().item()
124     return {
125         "loss": train_loss / len(train_ds),
126         "accuracy": 100.0 * correct / len(train_ds.dataset),
127     }
128
129
130 def test(model: nn.Module, test_ds: DataLoader, device: torch.device):
131     with torch.no_grad():
132         model.eval()
133         test_loss, correct = 0, 0
134         for data, target in tqdm(test_ds, desc="Test"):

```

```

135         data, target = data.to(device), target.to(device)
136         output = model(data)
137         test_loss += F.nll_loss(output, target).item()
138         pred = output.max(1, keepdim=True)[1]
139         correct += pred.eq(target.view_as(pred)).sum().item()
140     return {
141         "loss": test_loss / len(test_ds),
142         "accuracy": 100.0 * correct / len(test_ds.dataset),
143     }
144
145
146 def main(epochs: int = 2):
147     if not torch.cuda.is_available():
148         raise RuntimeError("CUDA is not available.")
149
150     device = torch.device("cuda")
151
152     batch_size = get_batch_size(
153         model=ResNet(),
154         device=device,
155         input_shape=IMAGE_SHAPE,
156         output_shape=(NUM_CLASSES,),
157         dataset_size=DATASET_SIZE,
158     )
159
160     train_ds, test_ds = get_datasets(batch_size=batch_size)
161     model = ResNet().to(device)
162
163     optimizer = optim.Adam(model.parameters(), lr=1e-3)
164
165     for epoch in range(1, epochs + 1):
166         print(f"\nEpoch {epoch}/{epochs}")
167         train_result = train(
168             model=model, optimizer=optimizer, train_ds=train_ds, device=device
169         )
170         test_result = test(model=model, test_ds=test_ds, device=device)
171         print(
172             f'Train loss: {train_result["loss"]:.04f}\t'
173             f'accuracy: {train_result["accuracy"]:.2f}%\n'
174             f'Test loss: {test_result["loss"]:.04f}\t'
175             f'accuracy: {test_result["accuracy"]}%\n'
176         )
177
178

```

```
179 if __name__ == "__main__":  
180     main()
```

find_batch_size_example.py hosted with ♥ by GitHub

[view raw](#)

Finally, here are some articles and papers on the topic of batch size and its effects on deep neural networks.

- [Epoch vs Batch Size vs Iterations](#)
- [Effect of batch size on training dynamics](#)
- [What's the Optimal Batch Size to Train a Neural Network?](#)
- [On Large-Batch Training for Deep Learning: Generalization Gap and Sharp Minima](#)
- [Don't Decay the Learning Rate, Increase the Batch Size](#)

Deep Learning

Pytorch

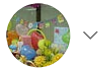
Batch Size

Programming

Open in app ↗



Search Medium



tds

Follow



Written by Bryan M. Li

55 Followers · Writer for Towards Data Science

Biomedical AI CDT student at the University of Edinburgh. bryanli.io

More from Bryan M. Li and Towards Data Science



Bryan M. Li in Towards Data Science

Accelerated TensorFlow model training on Intel Mac GPUs

A mini-guide on how to train TensorFlow model on MacBook Pro dGPU via TensorFlow PluggableDevice

5 min read · Oct 26, 2021



122



1





Matt Chapman in Towards Data Science

How I Stay Up to Date With the Latest AI Trends as a Full-Time Data Scientist

No, I don't just ask ChatGPT to tell me

★ · 8 min read · May 1

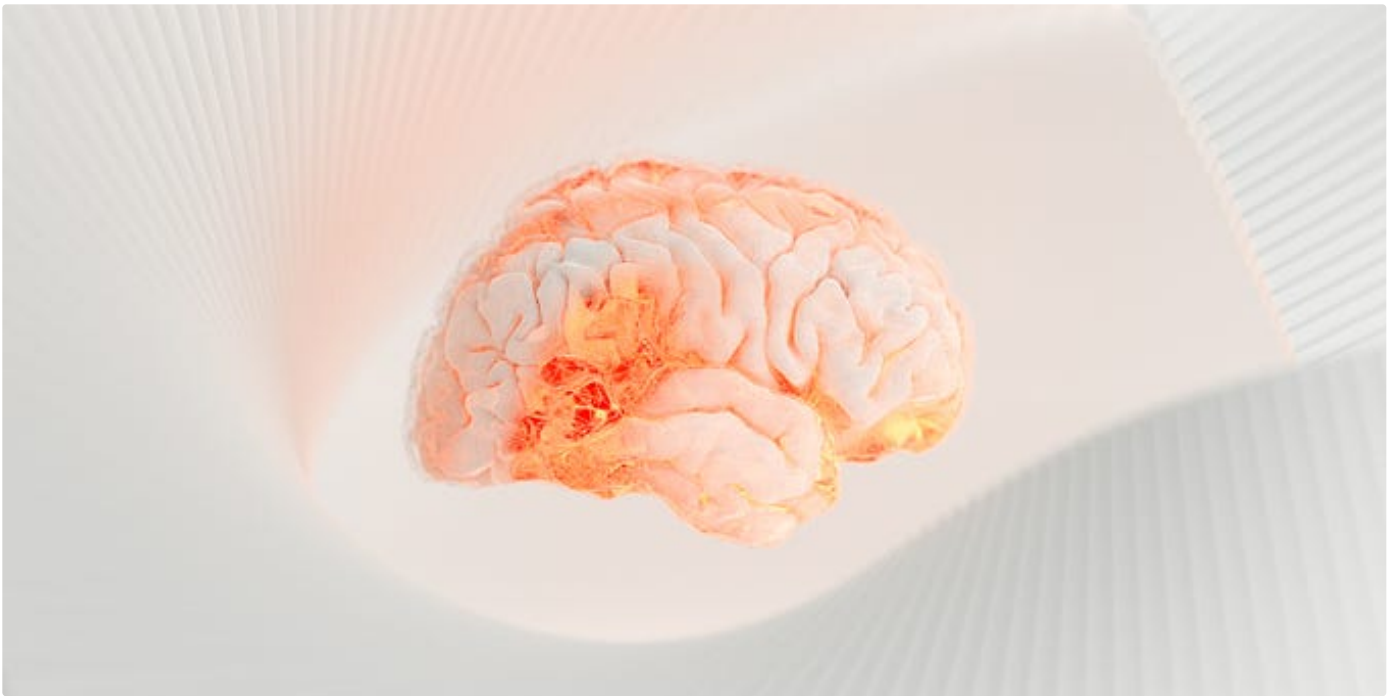


1.6K



21





Beatriz Stollnitz in Towards Data Science

How GPT Models Work

Learn the core concepts behind OpenAI's GPT models

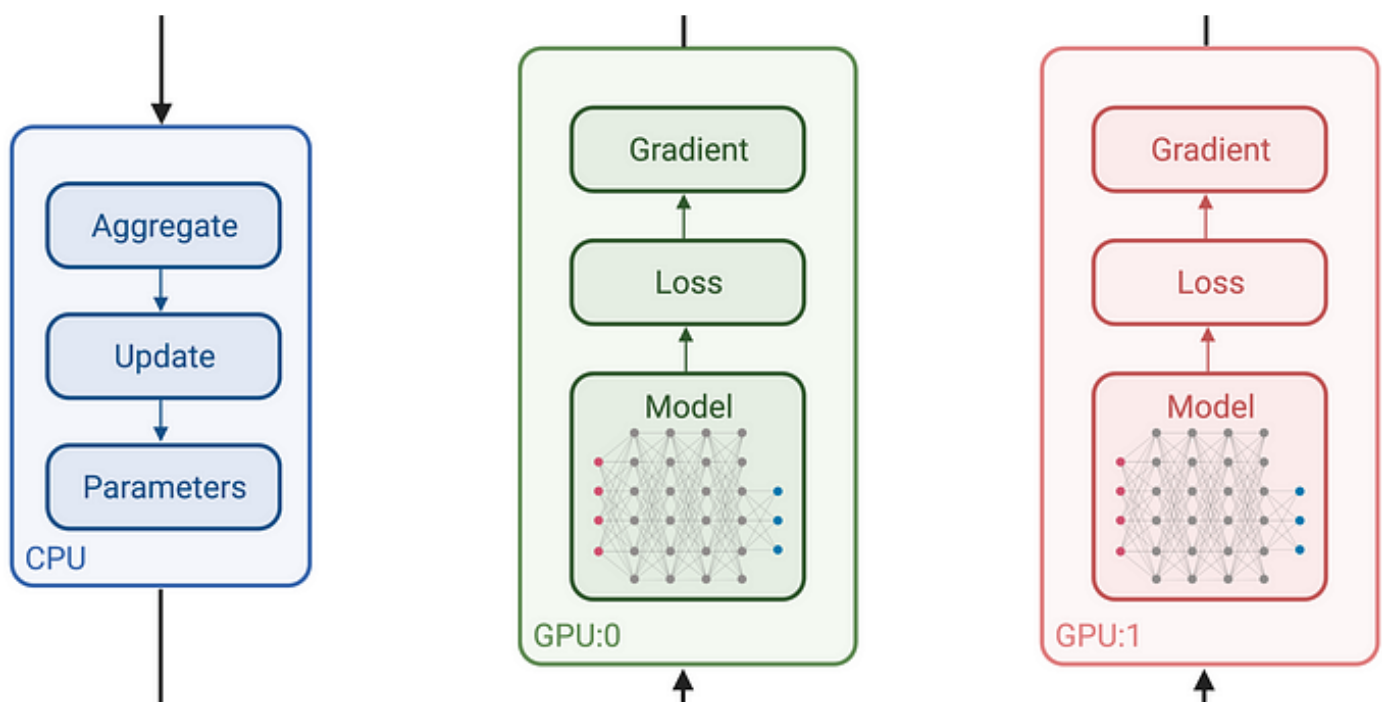
14 min read · May 20



1.1K



12





Bryan M. Li in Towards Data Science

Multi-GPUs and custom training loops in TensorFlow 2

A concise example of how to use `tf.distribute.MirroredStrategy` to train custom training loops model on multiple GPUs.

8 min read · Jun 1, 2021



84

[See all from Bryan M. Li](#)[See all from Towards Data Science](#)

Recommended from Medium





Cameron R. Wolfe in Towards Data Science

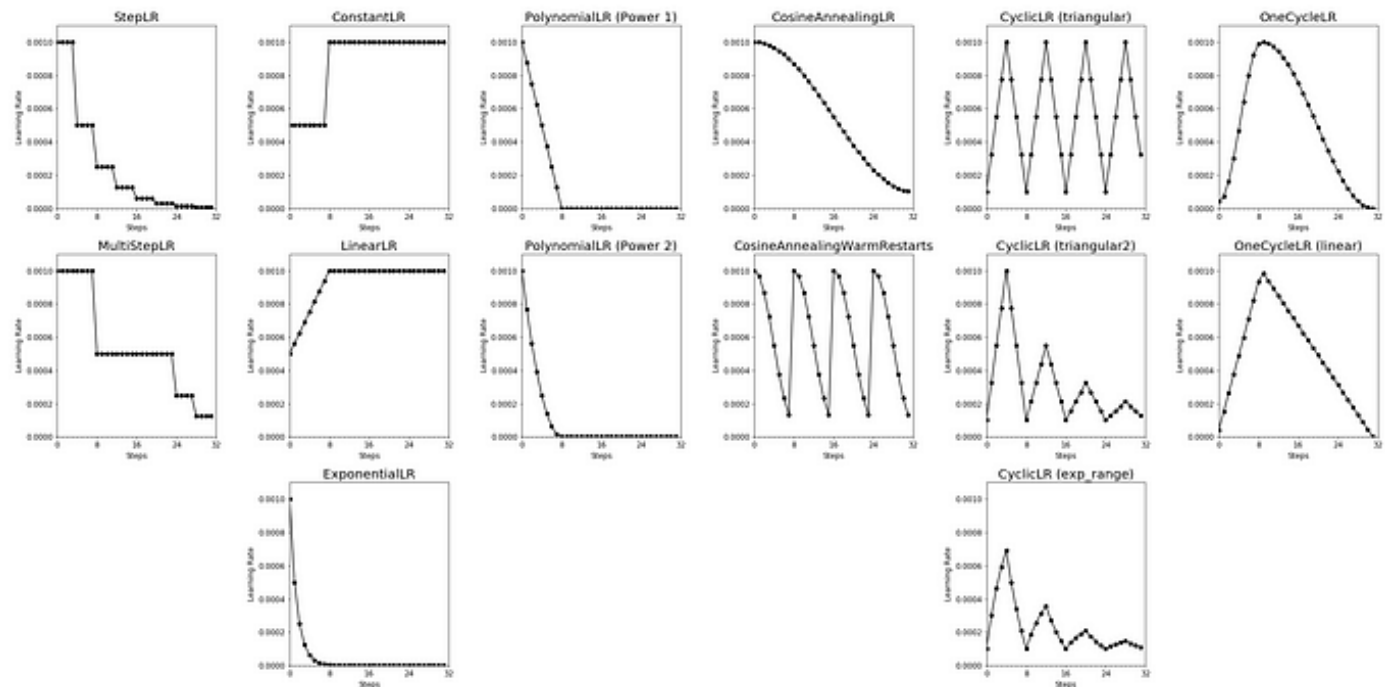
The Best Learning Rate Schedules

Practical and powerful tips for setting the learning rate

🌟 · 16 min read · Nov 16, 2022



99



Leonie Monigatti in Towards Data Science

A Visual Guide to Learning Rate Schedulers in PyTorch

LR decay and annealing strategies for Deep Learning in Python

🌟 · 9 min read · Dec 6, 2022



1K



6



Lists



Stories to Help You Grow as a Software Developer

19 stories · 85 saves



Leadership

30 stories · 31 saves



How to Run More Meaningful 1:1 Meetings

11 stories · 40 saves



Stories to Help You Level-Up at Work

19 stories · 68 saves



Thai Tran

Training PyTorch models on a Mac M1 and M2

PyTorch models on Apple Silicon M1 and M2

★ · 9 min read · Mar 25



17





Rukshan Pramoditha in Data Science 365

Determining the Right Batch Size for a Neural Network to Get Better and Faster Results

Guidelines for choosing the right batch size to maintain optimal training speed and accuracy while saving computer resources

🌟 · 4 min read · Sep 27, 2022



35



To redeem your free Amazon gift card click here!

Compute

Computation time on Intel Xeon 3rd Gen Scalable cpu: 0.024 s

SPAM

0.95

SPAM

0.00



Skanda Vivek in Towards Data Science

Transformer Models For Custom Text Classification Through Fine-Tuning

A tutorial on how to build a spam classifier (or any other classifier) by fine-tuning the DistilBERT model

🌟 · 4 min read · Jan 20

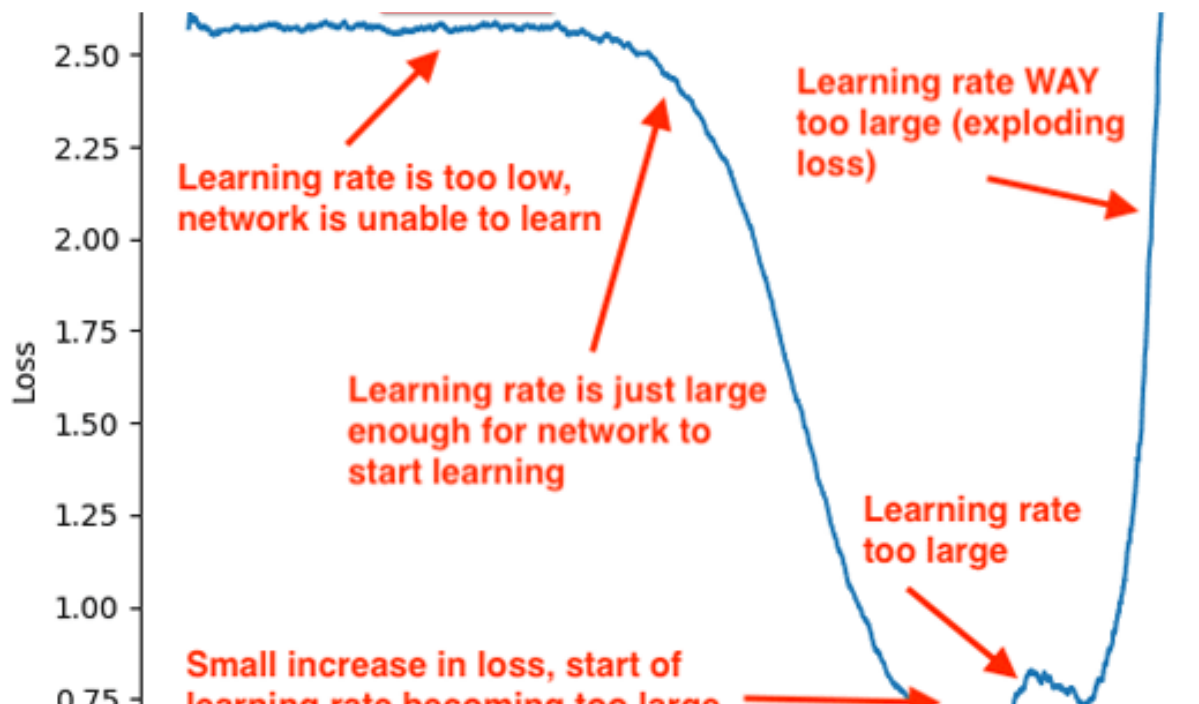


140



1





Moklesur Rahman

Learning Rate Scheduler in Keras

The learning rate is considered one of the most important hyperparameters for training deep learning models, but choosing it can be quite...

✦ · 7 min read · Jan 23



30



See more recommendations