



Checkpointing in ML Code

Dominik Viehhauser

Checkpoints

- Saving model weights during training

Checkpoints

- Saving model weights during training
- Save optimizer state if applicable (e.g. Adam)

Checkpoints

- Saving model weights during training
- Save optimizer state if applicable (e.g. Adam)
- General benefits:

Checkpoints

- Saving model weights during training
- Save optimizer state if applicable (e.g. Adam)
- General benefits:
 - Training can be resumed in case of failures

Checkpoints

- Saving model weights during training
- Save optimizer state if applicable (e.g. Adam)
- General benefits:
 - Training can be resumed in case of failures
 - Intermediate results can be compared and used (e.g. early stopping)

Checkpoints

- Saving model weights during training
- Save optimizer state if applicable (e.g. Adam)
- General benefits:
 - Training can be resumed in case of failures
 - Intermediate results can be compared and used (e.g. early stopping)
 - Be aware of storage requirements -> Avoid too frequent checkpointing

Benefits on the Cluster

- Increased resistance to failures: hardware, software, user error

Benefits on the Cluster

- Increased resistance to failures: hardware, software, user error
- Runtime limit of most jobs: up to 5 days

Benefits on the Cluster

- Increased resistance to failures: hardware, software, user error
- Runtime limit of most jobs: up to 5 days
 - Employ checkpointing and reduce max job time

Benefits on the Cluster

- Increased resistance to failures: hardware, software, user error
- Runtime limit of most jobs: up to 5 days
 - Employ checkpointing and reduce max job time
 - Splits a single long running job into multiple shorter ones with dependencies

Benefits on the Cluster

- Increased resistance to failures: hardware, software, user error
- Runtime limit of most jobs: up to 5 days
 - Employ checkpointing and reduce max job time
 - Splits a single long running job into multiple shorter ones with dependencies
 - Improves scheduling behavior of SLURM

Benefits on the Cluster

- Increased resistance to failures: hardware, software, user error
- Runtime limit of most jobs: up to 5 days
 - Employ checkpointing and reduce max job time
 - Splits a single long running job into multiple shorter ones with dependencies
 - Improves scheduling behavior of SLURM
 - Faster response time for updates/maintenance on the nodes

Benefits on the Cluster

- Increased resistance to failures: hardware, software, user error
- Runtime limit of most jobs: up to 5 days
 - Employ checkpointing and reduce max job time
 - Splits a single long running job into multiple shorter ones with dependencies
 - Improves scheduling behavior of SLURM
 - Faster response time for updates/maintenance on the nodes
 - Improve fairness of job scheduling

Checkpointing in PyTorch

Checkpointing in PyTorch

- PyTorch support saving and restoring model weights

```
1 import torch
2
3 model = ...
4
5 torch.save(model.state_dict(), filename)
6
7 model.load_state_dict(torch.load(filename))
```

Checkpointing in PyTorch

- PyTorch support saving and restoring model weights
- Multiple approaches possible

```
1 import torch
2
3 model = ...
4
5 torch.save(model.state_dict(), filename)
6
7 model.load_state_dict(torch.load(filename))
```

Checkpointing in PyTorch

- Store information about current epoch in filename

```
1 import torch
2
3 model = ...
4
5 if start_epoch > 0:
6     resume_epoch = start_epoch - 1
7     model.load_state_dict(torch.load("epoch-{} .pth".format(resume_epoch)))
8
9 for epoch in range(start_epoch, n_epochs):
10     ...
11     torch.save(model.state_dict(), "epoch-{} .pth".format(epoch))
```

Checkpointing in PyTorch

- Store information about current epoch in filename
- Also, possible to store multiple/all checkpoints in a single file

```
1 import torch
2
3 model = ...
4
5 if start_epoch > 0:
6     resume_epoch = start_epoch - 1
7     model.load_state_dict(torch.load("epoch-{:resume_epoch}.pth"))
8
9 for epoch in range(start_epoch, n_epochs):
10     ...
11     torch.save(model.state_dict(), "epoch-{:epoch}.pth")
```

Checkpointing in PyTorch

- Store information about current epoch in filename
- Also, possible to store multiple/all checkpoints in a single file
 - Easier automatic restart of last checkpoint

```
1 import torch
2
3 model = ...
4
5 if start_epoch > 0:
6     resume_epoch = start_epoch - 1
7     model.load_state_dict(torch.load("epoch-{:resume_epoch}.pth"))
8
9 for epoch in range(start_epoch, n_epochs):
10     ...
11     torch.save(model.state_dict(), "epoch-{:epoch}.pth")
```

Checkpointing in PyTorch

- Store information about current epoch in filename
- Also, possible to store multiple/all checkpoints in a single file
 - Easier automatic restart of last checkpoint
- Including a start epoch allows resuming training from a specific checkpoint

```
1 import torch
2
3 model = ...
4
5 if start_epoch > 0:
6     resume_epoch = start_epoch - 1
7     model.load_state_dict(torch.load("epoch-{} .pth".format(resume_epoch)))
8
9 for epoch in range(start_epoch, n_epochs):
10     ...
11     torch.save(model.state_dict(), "epoch-{} .pth".format(epoch))
```

Checkpointing in PyTorch

- Some optimizer also have states, e.g. Adam

```
1 import torch
2
3 model = ...
4 optimizer = ...
5
6 if start_epoch > 0:
7     resume_epoch = start_epoch - 1
8     model.load_state_dict(torch.load("epoch-{:resume_epoch}.pth")[ 'model' ])
9     optimizer.load_state_dict(torch.load("epoch-{:resume_epoch}.pth")[ 'optimizer' ])
10
11 for epoch in range(start_epoch, n_epochs):
12     ...
13     torch.save({
14         'optimizer': optimizer.state_dict(),
15         'model': model.state_dict(),
16     }, "epoch-{:epoch}.pth")
```

Links and sources

- Checkpointing in PyTorch https://pytorch.org/tutorials/beginner/saving_loading_models.html
- Checkpointing in Tensorflow <https://www.tensorflow.org/guide/checkpoint>

Thank you for your attention!