

quantized-federated-learning

December 20, 2023

0.1 Federated Learning Model with Quantization

```
[55]: import torch
from torchvision import models
import copy

# Load the saved model
model = models.resnet50(pretrained=False)
model.fc = torch.nn.Linear(model.fc.in_features, 6) # Change the output layer
↳ to match your task

# Load the model state_dict from the saved .pth file
model.load_state_dict(torch.load('updated_global_model_1.pth'))

# Set the model to evaluation mode
model.eval()

# Create a new instance of the model and load the state_dict to mimic clone
quantized_model = models.resnet50(pretrained=False)
quantized_model.fc = torch.nn.Linear(quantized_model.fc.in_features, 6) #
↳ Change the output layer to match your task
quantized_model.load_state_dict(copy.deepcopy(model.state_dict()))

# Set the cloned model to evaluation mode
quantized_model.eval()

# Quantize the weights of the cloned model to -1 or 1
for param in model.parameters():
    param.data = torch.sign(param.data)

# Calculate the number of parameters in the original and quantized models
quantized_model_num_params = sum(p.numel() for p in quantized_model.
↳ parameters())

torch.save(quantized_model.state_dict(), "binary_quantized_model.pth")
print(f"Quantized model size: {quantized_model_num_params} parameters")
```

```

float32_param_size = quantized_model_num_params * 4 # 4 bytes per float32
    ↪parameter
print(f"Memory occupied by float32 parameters: {float32_param_size} bytes")

int8_param_size = quantized_model_num_params * 1 # 1 byte per int8 parameter
print(f"Memory occupied by int8 parameters: {int8_param_size} bytes")

saved_percentage = ((float32_param_size - int8_param_size) /
    ↪float32_param_size) * 100
print(f"Memory saved by using int8 parameters: {saved_percentage:.2f}%")

```

Quantized model size: 23520326 parameters
Memory occupied by float32 parameters: 94081304 bytes
Memory occupied by int8 parameters: 23520326 bytes
Memory saved by using int8 parameters: 75.00%

```

[48]: import torch
from torchvision import datasets, transforms
from torch.utils.data import DataLoader
import torch.nn as nn
import torch.optim as optim
import time
from torchvision.models import resnet50, ResNet50_Weights

def test_model(data_dir, model_name):
    import torch
    from torchvision import datasets, transforms
    from torch.utils.data import DataLoader
    import torch.nn as nn

    # Define paths to your test dataset folder
    test_data_dir = data_dir # Update with your test dataset path

    # Define transformations for testing (similar to training)
    test_transforms = transforms.Compose([
        transforms.Resize(64),
        transforms.ToTensor()
    ])

    # Load the test dataset using ImageFolder with the defined transformations
    test_dataset = datasets.ImageFolder(root=test_data_dir,
    ↪transform=test_transforms)

    # Define the test dataloader

```

```

test_dataloader = DataLoader(test_dataset, batch_size=32, shuffle=False,
↪num_workers=4)

# Load the model architecture
model = resnet50(weights=None)
model.eval()

# Replace the final fully connected layer for transfer learning with the
↪same num_classes
num_fts = model.fc.in_features
num_classes = 6
model.fc = nn.Linear(num_fts, num_classes)

# Load the trained weights from the saved .pth file
model.load_state_dict(torch.load(model_name))
model.eval()

# Move the model to GPU if available
device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
model = model.to(device)

# Evaluate the model on the test dataset for both top-1 and top-3 accuracy
correct_top1 = 0
correct_top3 = 0

total = 0
with torch.no_grad():
    for images, labels in test_dataloader:
        images = images.to(device)
        labels = labels.to(device)

        outputs = model(images)
        _, preds = torch.topk(outputs, 3, dim=1) # Get top-3 predictions
        total += labels.size(0)
        for i in range(labels.size(0)):
            if labels[i] == preds[i, 0]: # Check top-1 accuracy
                correct_top1 += 1
            if labels[i] in preds[i]: # Check top-3 accuracy
                correct_top3 += 1

top1_accuracy = 100 * correct_top1 / total
top3_accuracy = 100 * correct_top3 / total
print(f'Top-1 Accuracy on the {data_dir} dataset: {top1_accuracy:.2f}%')
print(f'Top-3 Accuracy on the {data_dir} dataset: {top3_accuracy:.2f}%')

```

```
[49]: test_model('../Dataset/validation_data', "binary_quantized_model.pth")
```

Top-1 Accuracy on the ../Dataset/validation_data dataset: 99.50%
Top-3 Accuracy on the ../Dataset/validation_data dataset: 99.98%

```
[50]: test_model('../Dataset/test_data', "binary_quantized_model.pth")
```

Top-1 Accuracy on the ../Dataset/test_data dataset: 99.44%
Top-3 Accuracy on the ../Dataset/test_data dataset: 99.98%