



# **Module 2 - Session 1: Data and Model Building**

# Module 2 Overview

## What will we learn?

- Data Pipeline → handle large datasets
- Beyond Sequential → custom architectures
- Optimizers
- Device Management
- Building an image classifier



# Session 1: Data and Model Building

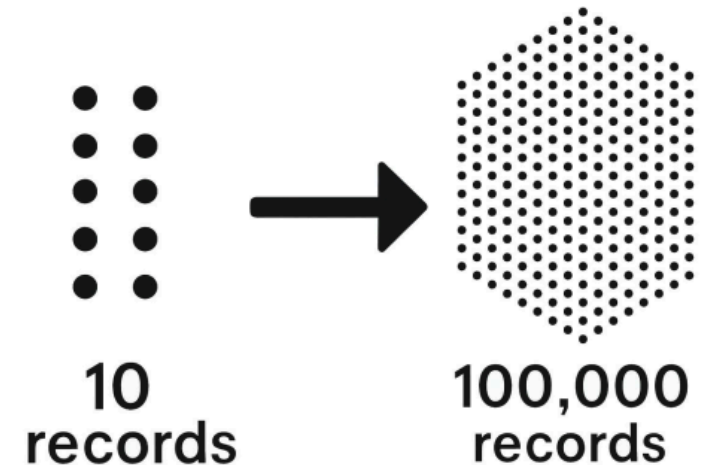
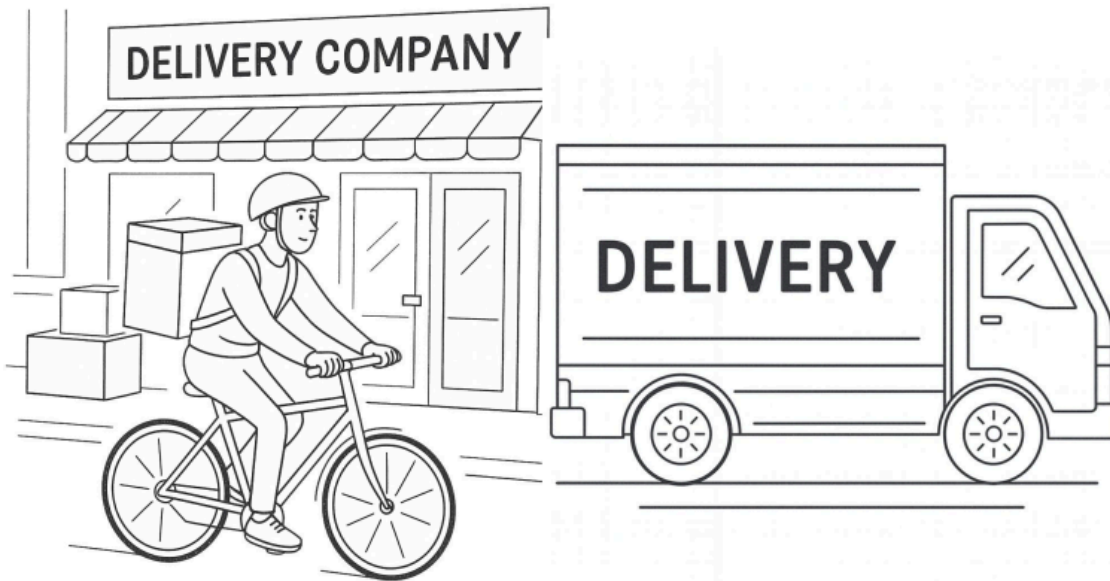
- Revisit the ML pipeline with a focus on PyTorch's data handling tools
- Learn about data management at scale
- Explore building custom model architectures beyond the Sequential API



# The Challenge: Large Datasets

**100,000 delivery records**

**Problem:** Loading all at once → runs out of memory



**Solution:** Work with data in batches

# PyTorch Data Utilities

## Three core tools:

1. **Transforms** - operations on each data point
2. **Dataset** - fetches samples from disk on demand
3. **DataLoader** - serves data in batches



# 1. Transforms

```
1 transform = transforms.Compose([  
2     transforms.ToTensor(),  
3     transforms.Normalize(mean=0.5, std=0.5)  
4 ])
```

**ToTensor:** converts to tensors, scales 0-255  $\rightarrow$  0-1

**Normalize:** centers around 0, scales using standard deviation



## 2. Dataset

```
1 dataset = MNIST(root='./data', train=True,  
2             download=True, transform=transform)
```



### Key features:

- Fetches samples from disk when asked
- Doesn't preload everything
- Handles where data lives, how to load samples, total count



# 3. DataLoader

```
1 dataloader = DataLoader(dataset, batch_size=64, shuffle=True)
```



**Batch size:** how many samples per batch

**Shuffle:** randomize order each epoch

**Makes training on large datasets possible**





# Complete Data Pipeline

```
1 # 1. Define transforms
2 transform = transforms.Compose([...])
3
4 # 2. Create dataset
5 train_dataset = MNIST(..., transform=transform)
6
7 # 3. Create dataloader
8 train_loader = DataLoader(train_dataset, batch_size=64, shuffle=True)
9
10 # 4. Use in training loop
11 for batch in train_loader:
12     images, labels = batch
13     # train model
```



# Beyond Sequential: Custom Models

## nn.Sequential

```
1 model = nn.Sequential(  
2     nn.Linear(1, 20),  
3     nn.ReLU(),  
4     nn.Linear(20, 1)  
5 )
```

## nn.Module

```
1 class MyModel(nn.Module):  
2     # defines layers  
3     def __init__(self):  
4         super().__init__()  
5         self.layer1 = nn.Linear(1, 20)  
6         self.layer2 = nn.Linear(20, 1)  
7  
8     # describes data flow  
9     def forward(self, x):  
10         x = self.layer1(x)  
11         x = F.relu(x)  
12         x = self.layer2(x)  
13         return x
```

**More control, same functionality**



# Calling the Model

**Don't call `model.forward()` directly**

**Do call `model(input)`**

PyTorch handles the forward call and essential bookkeeping



# Why `super().__init__()`?

## **Necessary for parameter tracking**

PyTorch needs to set up a system to track all learnable parameters (weights and biases)

Without it, PyTorch has nowhere to register your layers



# Training Loop Pattern

```
1 for batch in dataloader:
2     optimizer.zero_grad()      # Clear old calculations
3     outputs = model(inputs)    # Forward pass
4     loss = loss_fn(outputs, targets) # Measure error
5     loss.backward()            # Calculate gradients
6     optimizer.step()           # Update weights
```

**Order matters!** Don't swap these steps.



# Evaluation

```
1 model.eval() # Set to evaluation mode
2 with torch.no_grad(): # Disable gradient tracking
3     for batch in test_loader:
4         outputs = model(inputs)
5         # Calculate accuracy
```



## Two critical things:

- `model.eval()` - sets evaluation mode
- `torch.no_grad()` - disables gradient tracking



# Measuring Performance

## For classification: Accuracy

```
1 correct = (predictions == labels).sum().item()  
2 total = labels.size(0)  
3 accuracy = correct / total
```



**Count correct predictions / total predictions**



# To sum up

- Data pipeline: Dataset, DataLoader, Transforms
- Model building: `nn.Module`
- Training loop: `for batch in dataloader:`
- Evaluation: `model.eval()`, `torch.no_grad()`, accuracy





# What's Next?

In **Session 2: Loss Functions and Optimizers**, you learn:

- How loss functions measure error
- Cross-entropy loss for classification
- How optimizers use gradients to update weights
- Understanding backpropagation

