

Deep Learning: differentiable programming

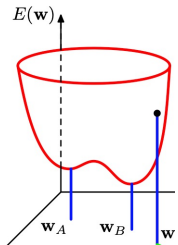
Christian Wolf

July 8th, 2021

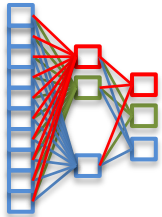
Content



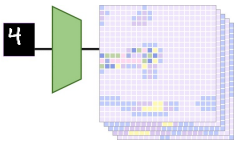
DL Frameworks and Tensors



Automatic differentiation

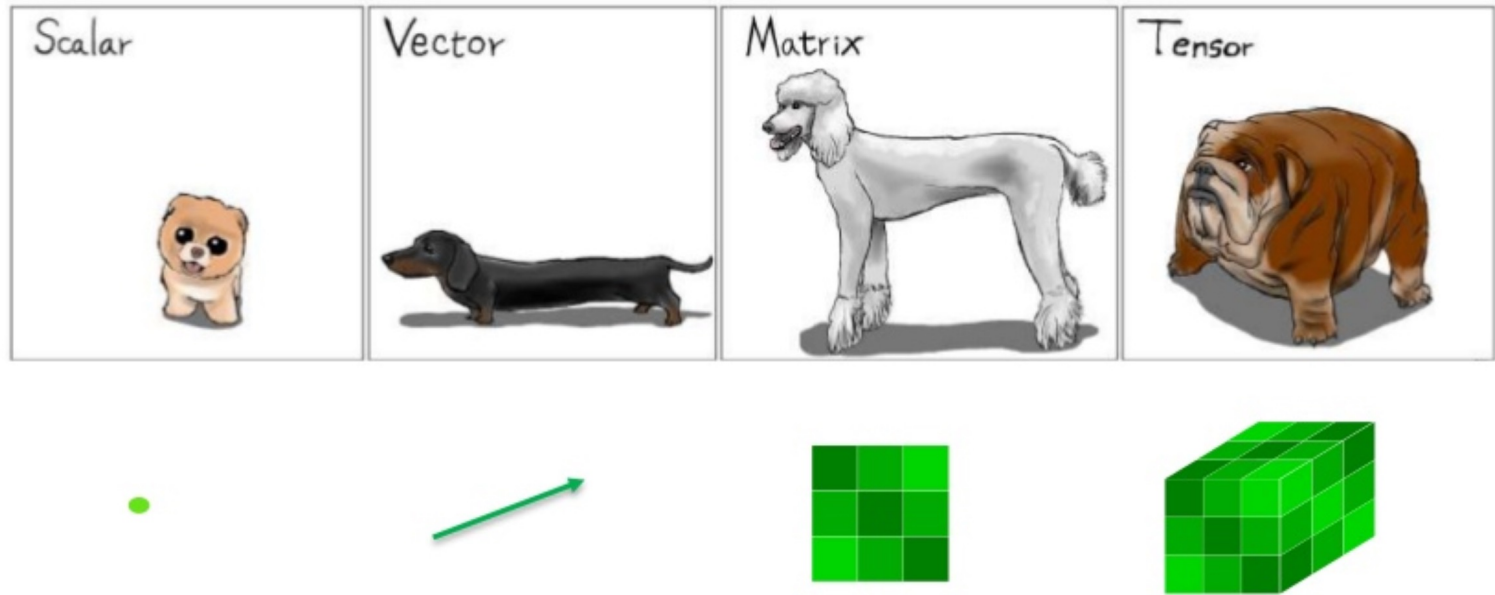


Example : MLP for MNIST



Exercise : CNN for MNIST

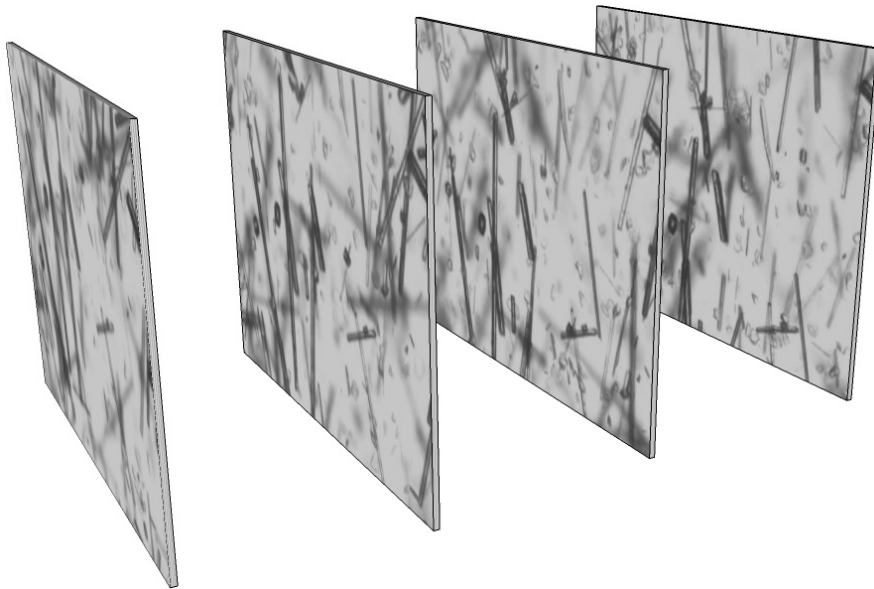
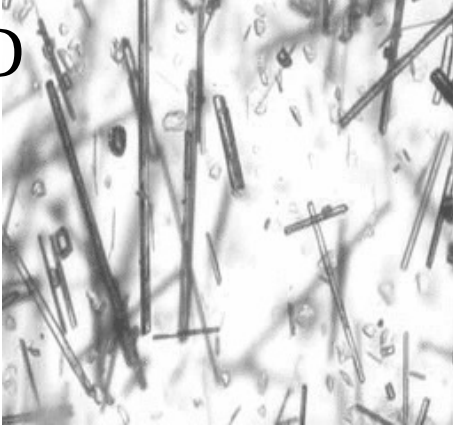
We manipulate tensors



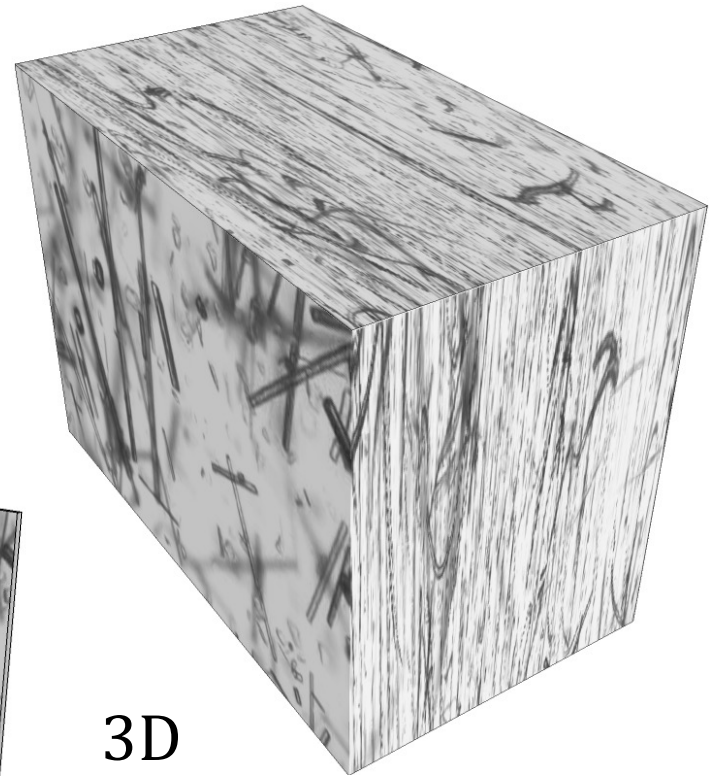
[Figure: Anima Anandkumar]

High dimensional tensors

2D



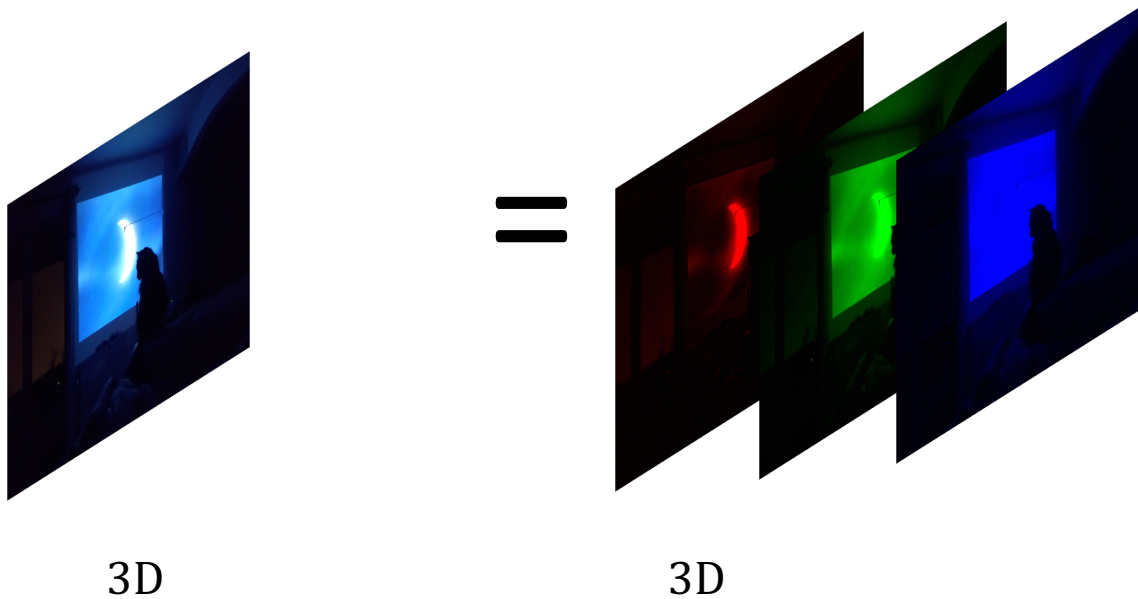
3D
(2D+t)



[Data: LAGEP]

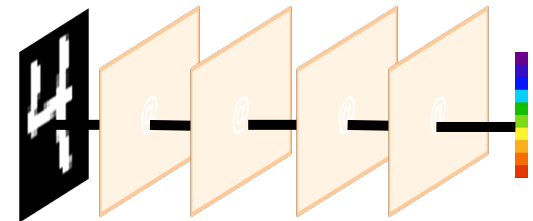
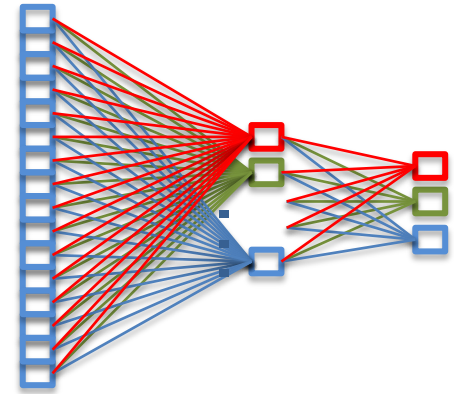
Images as tensors

A color image has 3 color channels (red, green, blue) and is therefore a 3D tensor.



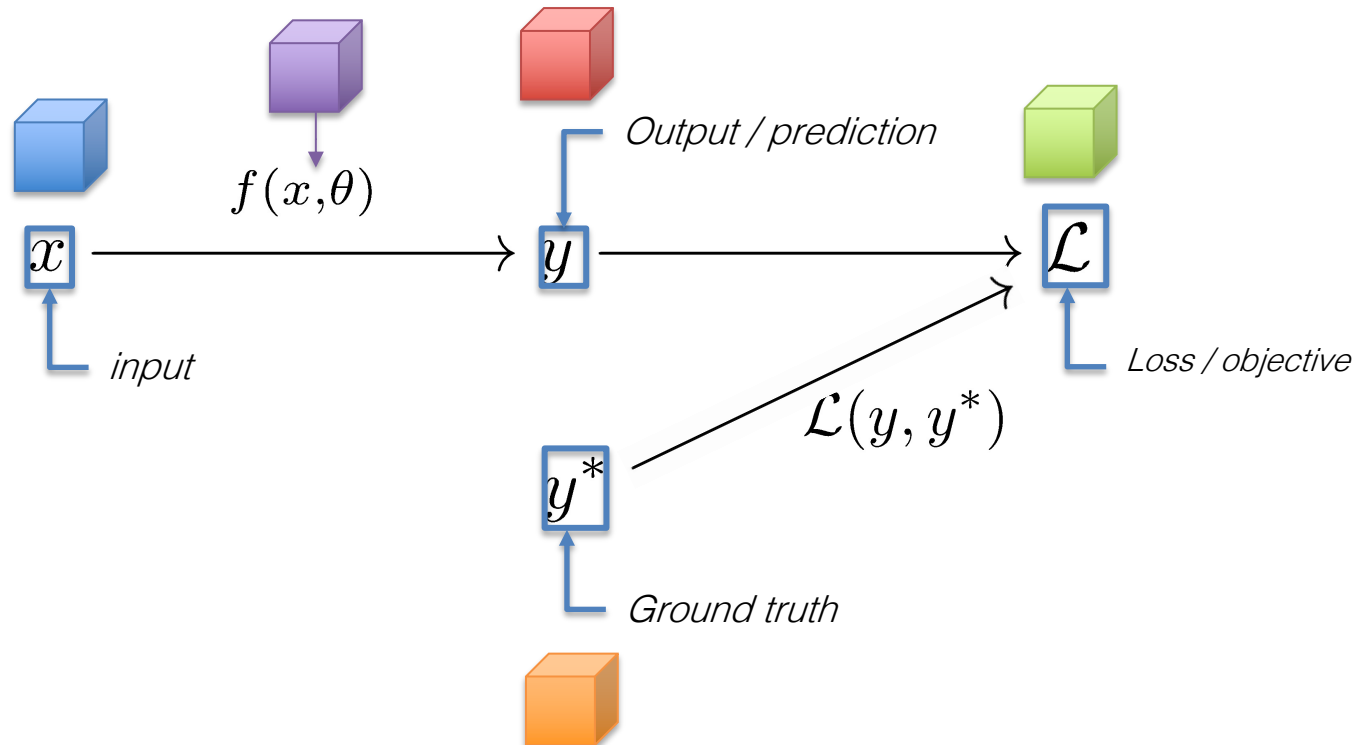
Tensors: examples

- Example of a tensor of dim 2 (input data, 1D signal)
 - Batch dimension (multiple samples)
 - Signal dimension
- Example of a tensor of dim 2 (output data, classification)
 - Batch dimension (multiple samples)
 - Prediction for different classes
- Example of a tensor of dim 3 (layer activation, 1D signal)
 - Batch dimension (multiple samples)
 - Signal dimension
 - Feature dimension
- Example of a tensor of dim 4 (layer activation, 2D image)
 - Batch dimension (multiple samples)
 - Spatial X dimension
 - Spatial Y dimension
 - Feature dimension
- Example of a tensor of dim 5 (input data, 2D+t video)
 - Batch dimension (multiple samples)
 - Spatial X dimension
 - Spatial Y dimension
 - Color channel dimension
 - Time dimension



Implementing a functional mapping

Inputs, outputs, layer activations, weights are tensors of different dimensions.



Main frameworks



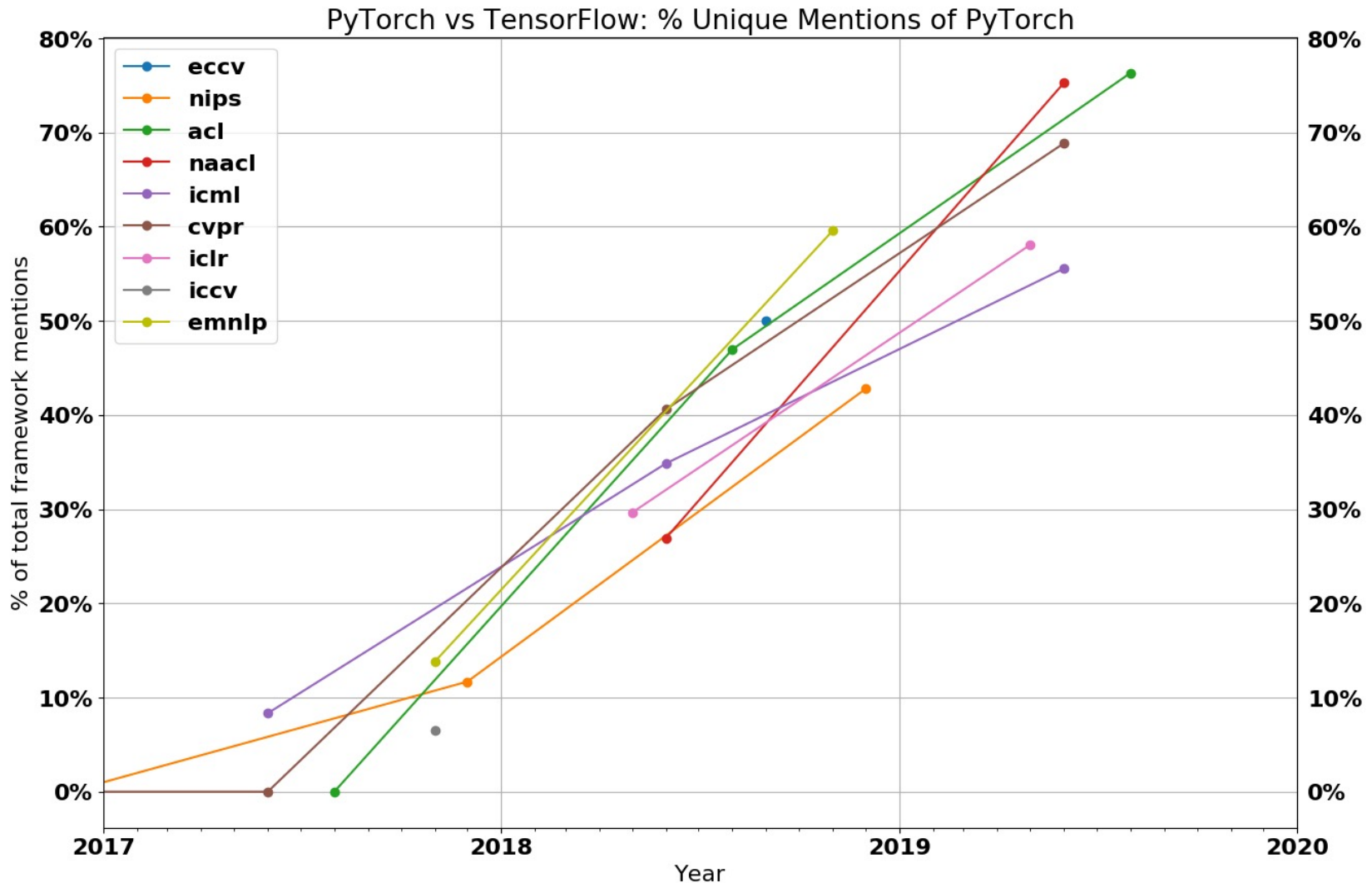
Tensorflow



PyTorch

- Both support execution and training on CPUs, GPUs, TPUs (google's machine learning hardware)
- Both use python.
- Tensorflow also supports C++, Swift.

PyTorch vs. Tensorflow



Creating tensors

```
1 # loading PyTorch
2 import torch
3
4 # create with given shape
5 torch.full((shape), value)
6 torch.full_like(other_tensor, value)
7
8 # create with given values
9 torch.tensor((values))
10 torch.tensor((values), dtype=torch.int16)
11
12 # create from numpy array
13 torch.from_numpy(numpyArray)
14
15 # Create zeros or ones
16 torch.zeros((shape))
17 torch.zeros_like(other_tensor)
18 torch.ones((shape))
19 torch.ones_like(other_tensor)
```

Creating tensors, tensor I/O

```
1 # Random tensors
2 torch.randn(3, 4)
3
4 # Tensor I/O
5 A = torch.load ("A.tensor")
6 torch.save (A, "A.tensor")
```

Tensor slicing

Slicing is similar to python (NumPy) Slicing or Matlab notation.
Example for a 2D tensor:

```
1 A[1,5]           # access an element (row, col)
2 A[:,5]           # column access
3 A[1,:]           # row access
4
5 A[1:6,:]         A[1:,:]       # range access (1:6 = 1,2,3,4,5)
6 A[:,0:-1]        # Negative index count backwards
7                  # -1 = last col/row
8
9 A==3             # provides a tensor of logical
10                # results
11
12 A[4:17,3] = B    # replace a slice
13
14 A[B==3]=4        # Set values in A to 4 at pos
15                # where there is a 3 in B
```

Manipulating tensors

```
1 # concatenate tensors
2 torch.cat((tensors), axis)
3
4 # split tensors into chunks of equal size
5 torch.split(tensor, splitSize, dim=0)
6
7 # reshape tensor w/o changing the data
8 torch.view(tensor, shape)
9
10 # Repeat along a given dimension
11 X.repeat(4,2)
12
13 # transpose tensor
14 torch.t(tensor) # 1D and 2D tensors
15 torch.transpose(tensor, dim0, dim1)
16
17 # Sorting
18 torch.sort(input, dim=-1)
```

Tensor math

```
1 # Overloaded operators
2 x = A+Y*Z-B          # * is elementwise mul
3
4 # Sum, product, min, max of all elements
5 torch.sum(tensor)    torch.min(tensor)
6 torch.prod(tensor)   torch.max(tensor)
7
8 # Linear algebra
9 torch.mm(A, B)        # Matrix multiplication
10 torch.inverse(tensor) # Matrix inversion
11 torch.det(tensor)     # Determinant
```

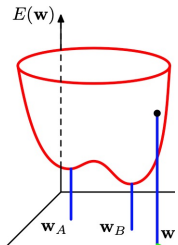
Elementwise operations

```
1 torch.exp(tensor)      torch.log(tensor)
2 torch.cos(tensor)      torch.cosh(tensor)
3 torch.sin(tensor)      torch.sinh(tensor)
4 torch.tan(tensor)      torch.tanh(tensor)
5
6 torch.add(tensor, tensor2) # or tensor+scalar
7 torch.div(tensor, tensor2) # or tensor/scalar
8 torch.mult(tensor, tensor2) # or tensor*scalar
9 torch.sub(tensor, tensor2) # or tensor-scalar
```

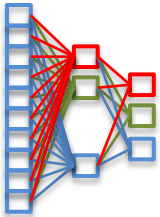
Content



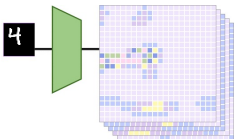
DL Frameworks and Tensors



Automatic differentiation



Example : MLP for MNIST



Exercise : CNN for MNIST

Gradient descent

One optimizer step:

$$\theta^{[t+1]} = \theta^{[t]} + \nu \nabla \mathcal{L}(h(x, \theta), y^*)$$

The gradient is a vector of partial derivatives:

$$\nabla \mathcal{L} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \theta_0} \\ \frac{\partial \mathcal{L}}{\partial \theta_1} \\ \vdots \\ \frac{\partial \mathcal{L}}{\partial \theta_N} \end{bmatrix}$$

Autograd

In PyTorch (and some other frameworks), Autograd performs automatic differentiation through a sequence of tensor instructions of an imperative language.

Let's consider a simple linear operation:

$$w = [5 \ 3], \quad x = [7 \ 2], \quad y = wx^T$$

The gradient of y w.r.t to x is given as

$$\nabla = \left[\frac{\partial y}{\partial x_i} \right] = \begin{bmatrix} 5 \\ 3 \end{bmatrix}$$

The gradient of y w.r.t to w is given as

$$\nabla = \left[\frac{\partial y}{\partial w_i} \right] = \begin{bmatrix} 7 \\ 1 \end{bmatrix}$$

Autograd

In PyTorch, we will first create the tensors:

```
1 w = torch.tensor([5, 3], dtype=float, requires_grad=True)
2 x = torch.tensor([7, 1], dtype=float, requires_grad=True)
```

The `requires_grad` flag ensures that all calculations are tracked. We perform the linear operation:

```
1 y = torch.dot(w, x)
```

Since the tensor y has been calculated as result of operations on tracked tensors, it has a gradient function:

```
1 print (y)
```

```
1 tensor(38., dtype=torch.float64, grad_fn=<DotBackward>)
```

Autograd

We now run a backward pass on the variable y , which calculates gradients w.r.t. to all involved tensors:

```
1 y.backward()
```

The gradients are attached to each variable:

```
1 print (x.grad)
2 print (w.grad)
```

```
1 tensor([5., 3.], dtype=torch.float64)
2 tensor([7., 1.], dtype=torch.float64)
```

Detaching tracking history

The tracking history uses memory in the tensor's space. If tracking is not used anymore for a tensor, it's tracking history can be detached:

```
1 print (y)
```

```
1 tensor(38., dtype=torch.float64, grad_fn=<DotBackward>)
```

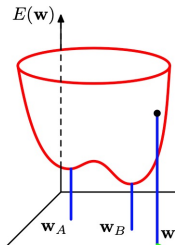
```
1 z = y.detach()  
2 print (z)
```

```
1 tensor(38., dtype=torch.float64)
```

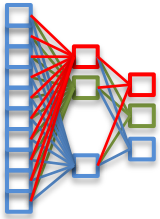
Content



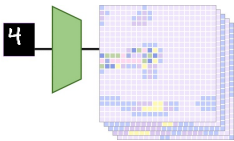
DL Frameworks and Tensors



Automatic differentiation



Example : MLP for MNIST

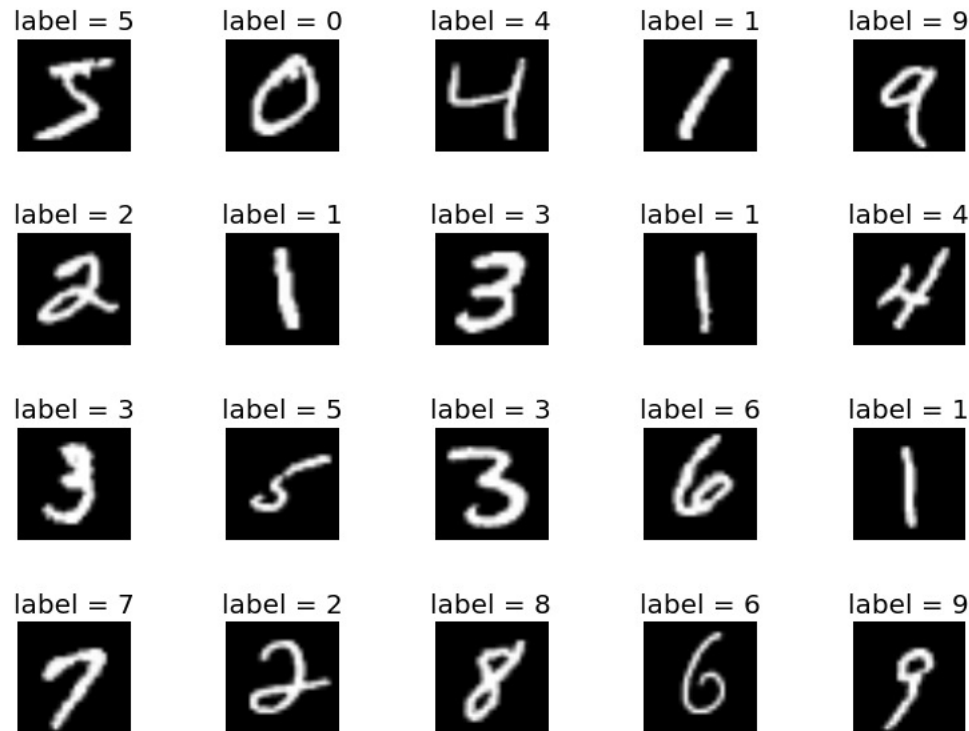


Exercise : CNN for MNIST

Example: the MNIST dataset

A dataset of handwritten digits introduced by Yann LeCun in 1999 with 60 000 training images and 10 000 test images.

One image is of size 28x28 pixels.



<http://yann.lecun.com/exdb/mnist/>

MNIST : MLP performance

linear classifier (1-layer NN)	none	12.0	LeCun et al. 1998
linear classifier (1-layer NN)	deskewing	8.4	LeCun et al. 1998
pairwise linear classifier	deskewing	7.6	LeCun et al. 1998

2-layer NN, 300 hidden units, mean square error	none	4.7	LeCun et al. 1998
2-layer NN, 300 HU, MSE, [distortions]	none	3.6	LeCun et al. 1998
2-layer NN, 300 HU	deskewing	1.6	LeCun et al. 1998
2-layer NN, 1000 hidden units	none	4.5	LeCun et al. 1998
2-layer NN, 1000 HU, [distortions]	none	3.8	LeCun et al. 1998
3-layer NN, 300+100 hidden units	none	3.05	LeCun et al. 1998
3-layer NN, 300+100 HU [distortions]	none	2.5	LeCun et al. 1998
3-layer NN, 500+150 hidden units	none	2.95	LeCun et al. 1998
3-layer NN, 500+150 HU [distortions]	none	2.45	LeCun et al. 1998

(Validation performance)

Writing Data Access

The MNIST dataset is supported by PyTorch through the Dataset class, which can even download the data from Yann Lecun's website:

```
1 train_dataset = datasets.MNIST("MNIST", True, download=True,
2   transform=transforms.Compose([
3   transforms.ToTensor(),
4   transforms.Normalize((0.1307,), (0.3081,))])
5   )
```

Writing Data Access

Let's recall training through gradient descent:

$$\theta^{[t+1]} = \theta^{[t]} + \nu \nabla \mathcal{L}(h(x, \theta), y^*)$$

The gradient is rarely (never?) taken over the whole dataset, but over a single sample, or batches (mini-batches) of a certain size. These batches are sample randomly from the dataset.

The actual shuffling and batching is performed by a built-in PyTorch DataLoader class, which uses an instance of our Dataset subclass:

```
1 train_dataset = datasets.MNIST("MNIST", True, download=True,
2   transform=transforms.Compose([
3   transforms.ToTensor(),
4   transforms.Normalize((0.1307,), (0.3081,))])
5 )
6
7 train_loader = torch.utils.data.DataLoader(train_dataset,
8   batch_size=BATCHSIZE, shuffle=True)
```

Tensor dimension conventions

PyTorch functions operate on multi-dimensional tensors and follow conventions on the order of dimensions.

- The first dimension is the batch dimension
 - ▶ Use 1 if you don't use batches (= batches of size 1).
 - ▶ Losses are reduced (sum or mean) over samples in a batch
- the second dimension is the channel dimension
 - ▶ Use 1 if you don't use channels (= single channels).
 - ▶ Channel arithmetics will be explained in detail in the section on convolutions.
- the following dimensions are application dependant, e.g. rows, columns in images.

The model

```
1 class LogRegression(torch.nn.Module):
2     def __init__(self):
3         super(LogRegression, self).__init__()
4
5         # input size to 10 output units
6         self.fc1 = torch.nn.Linear(28*28, 10)
7
8     def forward(self, x):
9         # Reshape from a 3D tensor (batchsize, 28, 28)
10        # to a flattened (batchsize, 28*28)
11        # 1 sample = 1 vector
12        x = x.view(-1, 28*28)
13        return self.fc1(x)
```

Set up the environment

```
1 # Instantiate the model
2 model = LogRegression()
3
4 # This criterion combines LogSoftMax and NLLLoss in
   one single class.
5 crossentropy = torch.nn.CrossEntropyLoss()
6
7 # Set up the optimizer: stochastic gradient descent
8 # with a learning rate of 0.01
9 optimizer = torch.optim.SGD(model.parameters(), lr
   =0.01)
10
11 # Init some statistics
12 running_loss = 0.0
13 running_correct = 0
14 running_count = 0
```

Iterative training

```
1 # Cycle through epochs
2 for epoch in range(100):
3
4     # Cycle through batches
5     for batch_idx, (data, labels) in enumerate(
6         train_loader):
7
8         # clear the gradients (they are accumulated)
9         optimizer.zero_grad()
10
11        # perform a forward pass and calc the loss
12        y = model(data)
13        loss = crossentropy(y, labels)
14
15        # backward pass: calculate the gradients
16        loss.backward()
17
18        # accumulate the loss, perform one SGD step
19        running_loss += loss.item()
20        optimizer.step()
```

Track training error

```
1  # Calculate the winner class
2  _, predicted = torch.max(y.data, 1)
3
4  # How many correct samples?
5  running_correct += (predicted == labels).sum().item()
6  running_count += BATCHSIZE
7
8  # Every 100 batches, print statistics
9  if (batch_idx % 100) == 0:
10
11     train_err = 100.0*(1.0 - running_correct / running_count)
12     print ('Epoch: %d batch: %5d ' % (epoch + 1, batch_idx +
13         1), end="")
14     print ('train-loss: %.3f train-err: %.3f' % (
15         running_loss / 100, train_err))
16     running_loss = 0.0
17     running_correct = 0.0
18     running_count=0.0
```

Logistic Regression: example output

```
1 Epoch: 1 batch:      1 train-loss: 0.026 train-err: 94.000
2 Epoch: 1 batch:    101 train-loss: 0.981 train-err: 25.020
3 Epoch: 1 batch:    201 train-loss: 0.538 train-err: 13.260
4 Epoch: 1 batch:    301 train-loss: 0.467 train-err: 12.700
5 Epoch: 1 batch:    401 train-loss: 0.425 train-err: 11.380
6 Epoch: 1 batch:    501 train-loss: 0.414 train-err: 11.820
7 Epoch: 1 batch:    601 train-loss: 0.384 train-err: 11.000
8 Epoch: 1 batch:    701 train-loss: 0.370 train-err: 10.180
9 Epoch: 1 batch:    801 train-loss: 0.374 train-err: 10.780
10 Epoch: 1 batch:    901 train-loss: 0.349 train-err: 9.480
11 Epoch: 1 batch:   1001 train-loss: 0.347 train-err: 9.840
12 Epoch: 1 batch:   1101 train-loss: 0.350 train-err: 9.920
13 Epoch: 2 batch:      1 train-loss: 0.347 train-err: 10.100
```

(...)

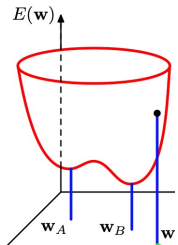
```
1 Epoch: 10 batch:   1001 train-loss: 0.277 train-err: 8.140
2 Epoch: 10 batch:   1101 train-loss: 0.277 train-err: 7.660
3 Epoch: 11 batch:      1 train-loss: 0.268 train-err: 7.700
```

This is training error, not validation error, i.e. **NOT** representative of the performance of the model!

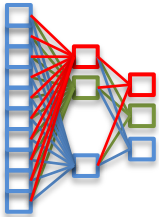
Content



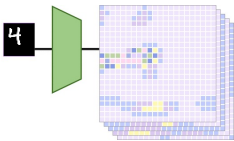
DL Frameworks and Tensors



Automatic differentiation



Example : MLP for MNIST



Exercise : CNN for MNIST

LeNet: example output

Exercise: create a convolution model similar to LeNet:

- One conv layer with 5×5 filters and 20 output channels.
- One conv layer with 5×5 filters and 50 output channels.
- One hidden layer with 500 units
- One output layer with 10 units (=digit classes)
- ReLU activation functions
- Max pooling (2×2) after each convolutional layer

Questions:

- How many filters does each layer have?
- How many trainable parameters does each layer have?
- What are the dimensions and sizes of the each output tensor?

LeNet: example output

```
1 Epoch: 1 batch:      1 train-loss: 0.024 train-err: 98.000
2 Epoch: 1 batch:    101 train-loss: 1.576 train-err: 36.140
3 Epoch: 1 batch:    201 train-loss: 0.747 train-err: 16.320
4 Epoch: 1 batch:    301 train-loss: 0.539 train-err: 12.880
5 Epoch: 1 batch:    401 train-loss: 0.481 train-err: 13.000
6 Epoch: 1 batch:    501 train-loss: 0.414 train-err: 11.280
7 Epoch: 1 batch:    601 train-loss: 0.386 train-err: 10.380
8 Epoch: 1 batch:    701 train-loss: 0.385 train-err: 10.900
9 Epoch: 1 batch:    801 train-loss: 0.363 train-err: 10.540
10 Epoch: 1 batch:    901 train-loss: 0.320 train-err:  9.120
11 Epoch: 1 batch:   1001 train-loss: 0.323 train-err:  8.920
12 Epoch: 1 batch:   1101 train-loss: 0.325 train-err:  9.400
13 Epoch: 2 batch:      1 train-loss: 0.304 train-err:  8.880
```

(...)

```
1 Epoch: 75 batch:    801 train-loss: 0.007 train-err:  0.000
2 Epoch: 75 batch:    901 train-loss: 0.007 train-err:  0.020
3 Epoch: 75 batch:   1001 train-loss: 0.008 train-err:  0.000
4 Epoch: 75 batch:   1101 train-loss: 0.009 train-err:  0.040
```

This is training error, not validation error, i.e. **NOT** representative of the performance of the model!