

PyTorch

Bryan



## PyTorch cheat sheet

```
import torch
```

```
from torch import nn
```

Almost everything in PyTorch is a module

```
from torch.utils.data import DataSet, DataLoader
```

Tensors:

```
Scalar = torch.tensor(7)
```

```
random_tensor = torch.rand(size=(3, 4))
```

Multipy 2 tensors : \*

torchvision cv lib → from torchvision import datasets, models, transforms  
torchtext text lib → from torchtext import datasets ...  
torch audio audio from torch audio ...  
torch rcc rcc engines

"cuda" > "mps" > "cpu"

device agnostic code

```
if torch.cuda.is_available():
    device = "cuda"
elif torch.backends.mps.is_available():
    device = "mps"
else:
    device = "cpu"
```

```
print(f"using device: {device}")
```

tensors have a device

```
x = torch.tensor(...)  
print(x.device)  
x = x.to(device)
```

```
torch.manual_seed(42)
```

```
torch.cuda.manual_seed(42)
```

## Neural Network Layers

```
torch.nn package  
from torch import nn
```

```
linear_layer = nn.Linear(m_features=10, out_features=10)
```

```
identity_layer = nn.Identity()
```

## Convolutional layers:

```
torch.nn.Convxd, x=1,2,3  
x = # of dimensions conv operates for
```

```
conv1d = nn.Conv1d(m_channels=1, out_channels=10, kernel_size=3)
```

## Transformer layers

```
transformer_model = nn.Transformer()  
dim embedding # of attention heads
```

```
transformer_encoder = nn.TransformerEncoderLayer(d_model=768, n_head=12)
```

```
transformer_decoder = nn.TransformerDecoderLayer(d_model=768, n_head=12)
```

```
transformer_decoder_stack = nn.TransformerDecoder(d_decoder_layer=transformer_decoder, num_layers=6)
```

## Recurrent Neural Networks :

### LSTMs & GRUs

`lstm_cell = nn.LSTMCell(input_size=10, hidden_size=10)`

`lstm_stack = nn.LSTM(input_size=10, hidden_size=10, num_layers=3)`

`gru_cell = nn.GRUCell(input_size=10, hidden_size=10)`

`gru_stack = nn.GRU(input_size=10, hidden_size=10, num_layers=3)`

### Activation Functions

`nn.ReLU`

`nn.Sigmoid`

`nn.Softmax`

`relu = nn.ReLU()`

### Loss functions

`nn.L1Loss`      `MAE`

`nn.MSELoss`      `MSE`

`nn.BCEWithLogitsLoss`      `Binary Cross Entropy`

`nn.CrossEntropyLoss`

## optimizers

torch.optim

```
torch.optim.SGD( lr = 0.1, params = model.parameters() )
```

```
torch.optim.Adam( lr = 0.001, params = model.parameters() )
```

Adam optimizer

model = nn...

optimizer = torch.optim...

80% 20% Split

using list slicing

## Creating a model

torch.nn.Module → forward() function

torch.nn.Sequential

```
from torch import nn
```

```
class LRM(nn.Module):
```

```
def __init__(self):
```

```
super().__init__()
```

```
self.linear_layer = nn.Linear(in_features=1, out_features=1)
```

```
def forward(self, x: torch.Tensor) -> torch.Tensor:
```

```
return self.linear_layer(x)
```

```
model = LRG()
```

```
model, model.state_dict()
```

using nn.Sequential

```
model = torch.nn.Sequential( nn.Linear(...))
```

model, model.state\_dict()

\* different state\_dict names

¶

Sequential gives default names

loss\_fn = nn.L1Loss()

```
optimizer = torch.optim.SGD(...)
```

epochs = 1000

```
x_train = x_train.to(device)
```

:

x-test...

y-train...

y-test...

```
model = model.to(device)
```

for epoch in epochs:

```
model.train()
```

y\_pred = model(x\_train)

loss = loss\_fn(y-pred, y-train)

optimizer.zero\_grad() ??

loss.backward()

optimizer.step()

Testing:

```
model1.eval()
```

```
with torch.inference_mode():
```

```
test_pred = model1(X_test)
```

```
test_loss = loss_fn(test_pred, y_test)
```

```
if epoch % 100 == 0:
```

```
print(f'Epoch: {epoch} | Train Loss: {train_loss} | Test Loss: {test_loss}')
```

$$\textcircled{1} \quad d\text{cost} / d\log \text{probs} \Rightarrow 32 \times 27$$

update log probs properly

$$\text{loss} = -\log \text{probs} (\text{range}(n) \cdot y_b] \cdot \text{mean}(c))$$

Averaging across batch size n

for the correct letter,  $d\text{loss}/\text{correct label idx} = -\frac{1}{n}$

The other gradients are zero

$$\textcircled{2} \quad d\text{cost} / d\text{probs} = \frac{d\text{cost}}{d\log\text{probs}} \cdot \frac{d\log\text{probs}}{d\text{probs}}$$

$$\log \text{probs} = \log(\text{probs}) \quad \text{vector to vector}$$

$$\frac{d\log\text{probs}}{d\text{probs}} = \left( \frac{1}{\text{probs}} \right)$$

$$y = \log(\pi)$$

2 operations

$$\textcircled{3} \quad d\text{cost} / d\text{counts\_sum\_inv}$$

Example:

$$\begin{aligned} C &= a * b & [ ] [ ] [ ] \\ & \left[ \begin{array}{|c|} \hline a_1 & a_2 & a_3 \\ \hline \end{array} \right] \times \left[ \begin{array}{|c|} \hline b_1 & b_2 & b_3 \\ \hline \end{array} \right] \\ a & \left[ \begin{array}{|c|} \hline 3 \times 3 \\ \hline \end{array} \right] \times b \left[ \begin{array}{|c|} \hline 3 \times 1 \\ \hline \end{array} \right] \\ & \left[ \begin{array}{|c|} \hline a_{11} \times b_1 & a_{12} \times b_1 & a_{13} \times b_1 \\ \hline \vdots & \vdots & \vdots \\ \hline a_{31} \times b_1 & a_{32} \times b_1 & a_{33} \times b_1 \\ \hline \end{array} \right] \end{aligned}$$

$$\frac{d\text{cost}}{d\log\text{probs}}$$

$$\frac{d\text{probs}}{d\log\text{probs}} \rightarrow \text{counts} *$$

$$d\text{counts\_sum\_inv}$$

$$\text{probs} = \frac{\text{vector}}{\text{Counts}} * \frac{\text{vector}}{\text{Counts\_sum\_inv}} \quad (32 \times 27) \quad (32 \times 1)$$

$$\text{Counts.shape} \quad 32 \times 27$$

$$\text{Counts\_sum\_inv.shape} \quad 32 \times 1$$

$$\textcircled{9} \quad \frac{1 \text{ cost}}{1 \text{ count}} \Rightarrow 32 \times 27$$

$$(32 \times 1) * (32 \times 27)$$

$$\frac{d \text{ poly}}{2 \text{ counts}} = \text{counts\_sum\_inv} \dots$$

column wise multiplication

$$\textcircled{5} \quad \frac{d \text{ cost}}{}$$

$$\frac{d \text{ counts\_sum}}{}$$

$$\xrightarrow{\frac{d \text{ counts\_sum\_inv}}{2 \text{ counts}}} -\frac{1}{( \text{counts\_sum})^2}$$

$$\textcircled{6}$$

$$\frac{d \text{ cost}}{d \text{ counts}} = \dots \frac{d \text{ counts\_sum}}{2 \text{ counts}}$$

$$\begin{array}{ccc} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{array} \rightarrow \left[ \begin{array}{cc} b_1 & - \\ - & b_2 \\ - & b_3 \end{array} \right]$$

$a$  replicate horizontally 27 times

$$32 \left( \begin{array}{c} | \\ \vdots \\ | \end{array} \right) \rightarrow 32 \left( \begin{array}{c} | \\ \vdots \\ | \end{array} \right)^{27}$$

$$(7) \frac{\partial \text{cost}}{\partial \text{norm\_logits}} = \dots$$

$$\frac{\partial \text{cost}}{\partial \text{norm\_logits}}$$

$$\text{Counts} = \text{norm\_logits} \cdot \exp(\cdot) \quad \text{elementwise}$$

$$\frac{\partial \text{counts}}{\partial \text{norm}} = \text{norm} \cdot \text{logits} \cdot \exp(\cdot)$$

$$(8) \frac{\partial \text{cost}}{\partial \text{logit\_maxes}} = \dots \frac{\partial \text{norm\_logits}}{\partial \text{logit\_maxes}}$$

$$\text{norm\_logits} = \text{logits} - \text{logit\_maxes}$$

$$(32, 27) \quad \text{logit\_maxes} = (32, 1)$$

$$\begin{bmatrix} C_{11} & C_{12} & C_{13} \\ C_{21} & C_{22} & C_{23} \\ C_{31} & C_{32} & C_{33} \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} - \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

copy

$$- \frac{\partial \dots}{\partial \dots}$$

$$(9) \frac{\partial \text{cost}}{\partial \text{logits}} = \frac{\partial \text{logit\_maxes}}{\partial \text{logits}} \rightarrow \text{need indices}$$

along columns

$\text{logit\_maxes} = \text{logits} \cdot \max(1, \dots) \text{ values} \Rightarrow \text{broadcasted to } 32 \times 27$

$32 \times 1 \quad 1 \text{ at } idx$

$\Sigma \text{ ad not cdx}$

$$\textcircled{10} \quad \frac{\partial C_{\text{cost}}}{\partial h} = \dots \underbrace{\frac{\partial \log_i h}{\partial h}}$$

$$\log_i h = h @ w + b,$$

$$\log_i h \quad (32, 27)$$

$$h \quad (32 \times 64)$$

$$w_2 \quad (64, 27)$$

$$b \quad (1 \times 27)$$

↑  
broadcasted column wise

$$d = a @ b + c$$

$$\begin{pmatrix} d_{11} & d_{12} \\ d_{21} & d_{22} \end{pmatrix} = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{pmatrix} + \begin{pmatrix} c_1 & c_2 \\ c_1 & c_2 \end{pmatrix}$$

$$d_{11} = a_{11} b_{11} + a_{12} b_{21} + c_1$$

$$d_{12} = a_{11} b_{12} + a_{12} b_{22} + c_2$$

$$d_{21} = a_{21} b_{11} + a_{22} b_{21} + c_1$$

$$d_{22} = a_{21} b_{12} + a_{22} b_{22} + c_2$$

$$\frac{\partial C}{\partial a_{11}} = \frac{\partial C}{\partial d_{11}} b_{11} + \frac{\partial C}{\partial d_{12}} b_{12}$$

$$\frac{\partial C}{\partial a_{12}} = \frac{\partial C}{\partial d_{11}} b_{21} + \frac{\partial C}{\partial d_{12}} b_{22}$$

$$\frac{\partial C}{\partial a_{21}} = \frac{\partial C}{\partial d_{21}} b_{11} + \frac{\partial C}{\partial d_{22}} b_{12}$$

$$\frac{\partial C}{\partial a_{22}} = \frac{\partial C}{\partial d_{21}} b_{21} + \frac{\partial C}{\partial d_{22}} b_{22}$$

$$\left( \begin{array}{cc} \frac{\partial C}{\partial a_{11}} & \frac{\partial C}{\partial a_{12}} \\ \frac{\partial C}{\partial a_{21}} & \frac{\partial C}{\partial a_{22}} \end{array} \right) \left( \begin{array}{cc} b_{11} & b_{21} \\ b_{12} & b_{22} \end{array} \right)$$

transposed



T

$$\frac{\partial C}{\partial d} \quad \textcircled{a} \quad b^T$$

order matters \*

$$\textcircled{11} \quad \frac{\partial C}{\partial b} = a^T \quad \textcircled{a} \quad \frac{\partial C}{\partial d}$$

$$\textcircled{12} \quad \frac{\partial C_{\text{tot}}}{\partial C} = \frac{\partial C_{\text{tot}}}{\partial d} \cdot \underset{\text{across columns}}{\text{Sum}(a)}$$

$$\textcircled{13} \quad \frac{\partial \text{cost}}{\partial \text{prc\_act}} = \dots \quad \frac{\text{act}}{\text{prc\_act}}$$

$$\text{act} = \tanh(\text{prc\_act})$$

$$\frac{\partial \text{act}}{\partial \text{prc\_act}} = 1 - (\text{prc\_act})^2$$

$$\textcircled{14} \quad \frac{\partial \text{cost}}{\partial \text{ln\_gain}} = \frac{\partial \text{prc\_act}}{\partial \text{ln\_gain}}$$

$\text{ln\_react}$  = ln<sub>base</sub> \* ln<sub>n\\_row</sub> + ln<sub>n</sub> bids

$$(32, 64) \quad (1, 64) \quad (32, 64) \quad (1, 64)$$

↑

lrc and costed

$$\text{dln\_raw} = ((1, 64) * (32 \times 64))$$

longer & direct

bn-mean

$$M_B = \frac{1}{m} \sum_{i=1}^m x_i$$

bn var

$$\sigma_B^2 = \frac{1}{m} \sum_{i=1}^m (x_i - M_B)^2$$

bndiff 2

bndraw

$$\frac{x_i}{\sigma_B^2} = \frac{(x_i - M_B)}{\sqrt{\sigma_B^2 + \epsilon}}$$

bndraw-inv (1, 64)

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv BN_{\gamma, \beta}(x_i)$$

Bessel's Correction

$$(32 \times 64) \quad (1, 64) \quad (32 \times 64)$$

(15)  $\text{dln diff} = \text{bn var\_inv} * \text{dln raw}$  ✓ w.r.t. bn\_raw

$$(1, 64) \quad (32, 64) \quad (32, 64)$$

(16)  $\text{dln\_var\_inv} = \text{bndiff} * \text{dln raw}$ , sum cor w.r.t. bn\_raw

$$(17) \frac{\partial \text{bnvar\_inv}}{\partial \text{bnvar}} = -0.5 \left( \text{bnvar} + (c-5)^{-1.5} \right) (1)$$

Bessel's correction: use for mini batches to get a better variance

(18)  $2 \text{ bn diff}_2$   $(32, 6^n)$

$2 \text{ bn var}$   $(1, 6^{n+1})$

$$\frac{2 \text{ bn-var}}{2 \text{ bndiff}}$$

$$\begin{matrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{matrix} = \begin{matrix} b_1 & b_2 \\ | & | \end{matrix}$$

$$b_1 = \left(\frac{1}{n-1}\right)(a_{11} + a_{21})$$

$$b_2 = \left(\frac{1}{n-1}\right)(a_{12} + a_{22})$$

$$\frac{2 \text{ bn-var}}{2 \text{ bn-diff}_2} = \left(\frac{1.0}{n-1}\right) (\text{torch. ones\_list (bn diff}_2))$$

Sum in forward pass: broad casting in the backward pass  
 Broadcasting in forward pass: sum in backward pass  
 ↗ replication  
 ↙ Variable reuse

(19) bndiff (again)

$$\frac{2 \text{ bn diff}_2}{2 \text{ bn diff}} = 2 \text{ bndiff}$$

(20)  $\frac{2 \text{ bndiff}}{2 \text{ bn var}} (1, 6^n)$  Sum in backward pass  $= -\text{torch. ones\_list (bndiff)} * \text{dbndiff}$

$$\frac{2 \text{ bndiff}}{2 \text{ bndiff}} = \text{dbndiff}$$

(22)  $\frac{1}{n} \sum h_{prob} \cdot \text{sum}(\partial, \text{keep...})$

↓  
broadcasting

$$\frac{\partial h_{prob}}{\partial h_{mean}} = \text{torch.zeros}_{} - \text{like}(h_{prob}) \left( \frac{1}{n} \right)$$

Repeat again ...

(23)  $\frac{\partial h_{prob}}{\partial emb} \Rightarrow$

(23)  $\frac{\partial h_{prob}}{\partial w_i} \Rightarrow$

(24)  $\frac{\partial h_{prob}}{\partial b_i} \Rightarrow$

(25)  $\frac{\partial emb}{\partial emb} \stackrel{\text{unflatten}}{\Rightarrow} \text{demb.cat().view(emb.shape)}$

(26) ! plucking out rows

using indexing

Gradients for 2+ occurrences must add

$$tC = \text{torch.zeros}_{} - \text{like}(C)$$

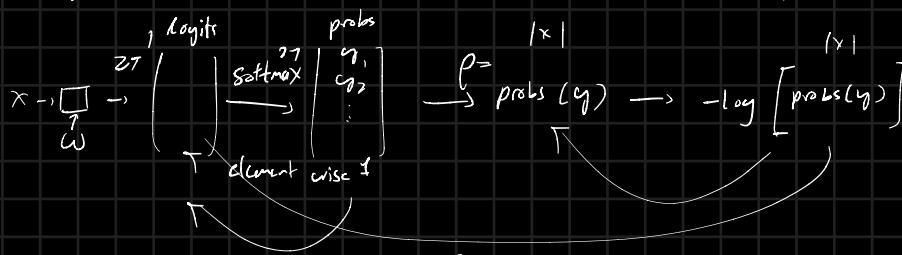
for k in range(C.size().shape[0]):

for j, n in enumerate(C.size().shape[1]):

$$inx = x_B(k, j)$$

$$dC[inx] += \text{demb}[k, j]$$

All together :



Cost

$|x|$

$|y|$

$$p = \frac{1}{\sum_i p_i}$$

$$\text{prob}(y) \rightarrow -\log \left[ \text{prob}(y) \right]$$

$$\text{loss} = -\log p_y$$

$$= -\log \left( \frac{e^{l_y}}{\sum_i e^{l_i}} \right)$$

$\left( \begin{array}{c} l_1 \\ l_2 \\ \vdots \\ l_n \end{array} \right)$   $\left( \begin{array}{c} 0 \\ 1 \\ 0 \\ 0 \end{array} \right)$   $\text{ith idx } \left( \begin{array}{c} 1 \\ 0 \\ \vdots \\ 0 \end{array} \right)$

$$\frac{\partial C}{\partial x} = \frac{1}{x}$$

$$\frac{\partial C}{\partial l_i} = \frac{\partial}{\partial l_i} \left( -\log \left( \frac{e^{l_y}}{\sum_j e^{l_j}} \right) \right)$$

$$\left( \begin{array}{c} 0 \\ -\frac{1}{\sum_j e^{l_j}} \\ 0 \\ \vdots \\ 0 \end{array} \right) \quad (\text{derivative of softmax})$$

$$\left( \begin{array}{c} 1 \\ 0 \\ \vdots \\ 0 \\ 0 \end{array} \right)$$

$$\frac{\partial C}{\partial \text{prob}} \quad \frac{\partial \text{prob}}{\partial \text{logits}}$$

$$p^{\text{prob}}_i = \frac{e^{l_i}}{\sum_j e^{l_j}}$$

$$\text{for all } i, \frac{\partial p^{\text{prob}}_i}{\partial \text{logits}} = \left( \frac{1}{\sum_j e^{l_j}} \right) \underbrace{\left( e^{l_i} \right)}_{\text{constant}}$$

$$\frac{\partial C}{\partial \text{logits}} \rightarrow \left( \frac{\partial C}{\partial \text{prob}} \right) \left( \frac{\partial \text{prob}}{\partial \text{logits}} \right)$$

$$p^{\text{prob}}_i = 1$$

$$i=1 \quad \frac{e^{l_1}}{(e^{l_1} + e^{l_2} + \dots + e^{l_n})}$$

$$\frac{\partial C}{\partial \theta_i} = \frac{2}{2\theta_i} \left( -\log \left( \frac{e^{\theta_i}}{\sum_j e^{\theta_j}} \right) \right)$$

$$= - \frac{\sum_j e^{\theta_j}}{e^{\theta_i}} \frac{1}{\theta_i} \left( \frac{e^{\theta_i}}{\sum_j e^{\theta_j}} \right)$$

$i = y$  :

$$= \frac{\sum_j e^{\theta_j}}{e^{\theta_y}} \left[ 0 \times \frac{1}{\sum_j e^{\theta_j}} + -\frac{e^{\theta_y}}{\left( \sum_j e^{\theta_j} \right)^2} \right]$$

$$= \frac{e^{\theta_i}}{\sum_j e^{\theta_j}} = p_i$$

$i = y$

$$= \frac{\sum_j e^{\theta_j}}{e^{\theta_y}} \left( \frac{e^{\theta_y}}{\sum_j e^{\theta_j}} + -\frac{e^{\theta_y} e^{\theta_i}}{\left( \sum_j e^{\theta_j} \right)^2} \right)$$

$$= p_i - 1$$

For batches, average all examples, and divide by n

Batch normalization:

Calculate  $\frac{\partial L}{\partial \theta}$  proba, given  $\frac{\partial L}{\partial \theta}$  preact

$$\mu = \frac{1}{m} \sum x_i$$

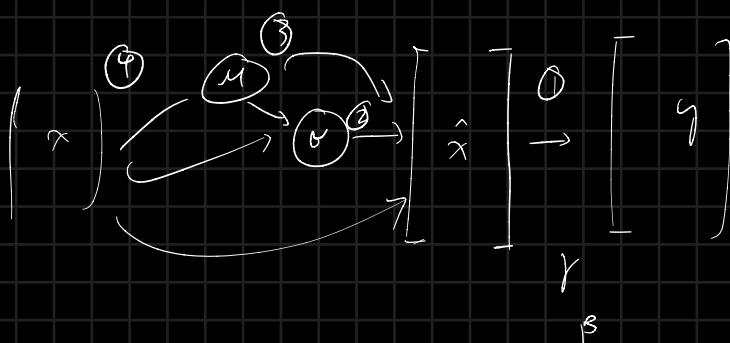
$$\sigma^2 = \frac{1}{m-1} \sum (x_i - \mu)^2$$

$$\hat{x}_i = \frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$y_i = \gamma \hat{x}_i + \beta$$

Having  $\frac{\partial L}{\partial y_i}$

$$\frac{\partial L}{\partial \hat{x}_i} = \frac{\partial L}{\partial \hat{x}_i} \frac{\partial \hat{x}_i}{\partial x_i}$$



$$\textcircled{1} \quad \frac{\partial y}{\partial \hat{x}} = \gamma$$

$$\textcircled{2} \quad \sigma^2 = \frac{1}{m-1} \sum (x_i - \mu)^2$$

$$\textcircled{3} \quad \left( \frac{\partial \sigma}{\partial x_i} \right) (\sigma) =$$

$$\textcircled{2} \quad \frac{\partial L}{\partial \sigma^2} = \sum_i \frac{\partial L}{\partial \hat{x}_i} \frac{\partial \hat{x}_i}{\partial \sigma^2}$$

$$= \gamma \cdot \frac{\partial L}{\partial \hat{x}_i} \frac{\partial \hat{x}_i}{\partial \sigma^2} \left( (x_i - \mu)(\sigma^2 + \epsilon)^{-\frac{1}{2}} \right)$$

$$= -\frac{1}{2} \gamma \frac{\partial L}{\partial \hat{x}_i} \cdot (x - \mu)(\sigma^2 + \epsilon)^{-\frac{3}{2}}$$

$$\textcircled{3} \quad \frac{\partial L}{\partial \mu} = \sum \left( \frac{\partial L}{\partial \hat{x}_i} \right) \frac{\partial \hat{x}_i}{\partial \mu} + \left( \frac{\partial L}{\partial \sigma^2} \right) \frac{\partial \sigma^2}{\partial \mu}$$

$$\frac{\partial \hat{x}_i}{\partial \mu} = \frac{1}{2m} \left( \frac{\gamma}{\sqrt{\sigma^2 + \epsilon}} \right) = -(\sigma^2 + \epsilon)^{-\frac{1}{2}}$$

$$\frac{\partial \sigma^2}{\partial \mu} = \frac{2}{2m} \left( \frac{1}{m-1} \sum_{i=1}^m (x_i - \mu)^2 \right)$$

$$= \frac{2}{m-1} \sum (x_i - \mu)$$

$$= \frac{2}{m-1} \left( \sum x_i - \sum \mu \right)$$

$$= \left( \frac{2}{m-1} \right) \left( m\mu - m\mu \right)$$

$$= \underline{\underline{\delta}}$$

$$\frac{\partial L}{\partial \mu} = - \sum \frac{\partial L}{\partial \hat{x}_i} \cdot \gamma \cdot (\sigma^2 + \epsilon)^{\frac{1}{2}}$$

4

 $x \rightarrow x_i \rightarrow 32$  arrowsexactly parallel  
and  $\frac{\partial x_i}{\partial x} = 1$ 

$$\frac{\partial c}{\partial x_i} = \frac{\partial c}{\partial x_i} \frac{\partial \hat{x}_i}{\partial x_i} + \frac{\partial c}{\partial u} \frac{\partial u}{\partial x_i} + \frac{\partial c}{\partial \alpha^2} \frac{\partial \alpha^2}{\partial x_i}$$

$$\frac{\partial \hat{x}_i}{\partial x_i} = \frac{d}{dx_i} \left( \frac{x_i - \mu}{\sqrt{\alpha^2 + \epsilon}} \right) = (\alpha^2 + \epsilon)^{-\frac{1}{2}}, \quad \frac{du}{dx_i} = \frac{d}{dx_i} \left( \frac{1}{m} \sum_j x_j \right) = \frac{1}{m}$$

$$\frac{\partial \alpha^2}{\partial x_i} = \frac{2}{\partial x_i} \left( \frac{1}{m-1} \sum_j (x_j - \mu)^2 \right) = \left( \frac{2}{m-1} x_i - \mu \right)$$

$$\left( \frac{\partial c}{\partial \hat{x}_i} \right) \left( \alpha^2 + \epsilon \right)^{-\frac{1}{2}} + \left( \sum_{j \neq i} \frac{\partial c}{\partial x_j} \cdot \gamma \cdot (\alpha^2 + \epsilon)^{-\frac{1}{2}} \right) \cdot \frac{1}{m}$$

$$+ \left( -\frac{1}{2} \gamma \sum_j \frac{\partial c}{\partial x_j} \cdot (x_j - \mu) (\alpha^2 + \epsilon)^{-\frac{3}{2}} \right) \left( \frac{2}{m-1} (x_i - \mu) \right)$$



torch vision => computer vision      torch. float Tensor(np.array(>))

class LinearModel(nn.Module):

def \_\_init\_\_(self, input\_size, output\_size):

super(LinearModel, self).\_\_init\_\_()

self.w = nn.Parameter(torch.randn(output\_size, input\_size))

def forward(self, x):

return x @ self.w.t()

output.detach() (removes tensor class)

① linear = nn.Linear(inputs, outputs)

linear.weight.data

output = linear(data)

② loss\_function = nn.MSELoss()

③ optimizer = torch.optim.SGD(linear.parameters(), lr=0.01)

loss.backward() (calculates gradients)

optimizer.step() (Takes a learning step)

Forward pass

(calculate loss)

reset stored gradients to zero

loss.backward()

take learning step with optimizer.step()

loss function expects 2 arguments of some form

loss function found in torch.nn

class Model(nn.Module):

def \_\_init\_\_(self, ...):

super(Model, self).\_\_init\_\_()

self.linear1 = nn.Linear(input\_size, hidden\_size)

self.linear2 = nn.Linear(hidden\_size, output\_size)

def forward(self, x):

self.h1 = self.linear1(x)

self.h2 = torch.Sigmoid(h1)

self.h3 = self.linear2(self.h2)

return self.h3

logistic = Model(input\_size=2, output\_size=1, hidden\_size=2)

loss\_fn = nn.BCEWithLogitsLoss()

optimizer = torch.optim.Adam(logistic.parameters(), lr=1e-2)

Class SmDataset(Dataset):

def \_\_init\_\_(self, num\_data):

self.x\_data = ...

self.y\_data = ...

def \_\_getitem\_\_(self, index):

''' Get a new batch, passes in the index

''' return self.x\_data[index], self.y\_data[index]

def \_\_len\_\_(self):

return self.x\_data.shape[0]

Torch dataset + dataloaders

Data loader object

-- init --

--getitem--

--len--

Data / DataLoaders

torch.utils.data import DataLoader

torch.utils.data.Dataset import Dataset

torchvision.datasets import MNIST

# Convolutional Neural Networks

PIL - Python Imaging Library

torch.utils.data.random\_split([total], valid), generator = torch.Generator().manual\_seed(42))

Seed is important when saving!

LocNet5

```
transform = transforms.Compose([
    transforms.Resize(32),
    transforms.ToTensor(),
    transforms.Normalize(...)])
```

nn.Conv2D(channels\_in, out\_channels, kernel\_size)

(batch, in\_channels, height, width)

(batch, out\_channels, out\_height, out\_width)

model.train()

model.eval()

=> Disables Dropout, Batch Norm (running mean & std dev)

fx.argmax(1).sumc().item()

PIL:

Image.open(...).convert('RGB')

# of filters

(out\_channels, in\_channels, k\_height, k\_width)

(batch, in\_channels, im\_height, im\_width)

nn.Conv2D(in\_channels, out\_channels, kernel\_size, stride, padding, bias=True)

# Transfer Learning

Weight freezing : Freeze pretrained model from being updated in backward pass

torch.vision.models as models

num\\_ftrs = resnet.fc.in\\_features

resnet.fc = nn.Linear(num\_ftrs, 10).to(device)

torch.load("model.pth")

model = torch.load(...)  
model[key]

model.load\_state\_dict(torch.load(save\_path[variable]))  
optimizer.load\_state\_dict(torch.load(save\_path[variable]))

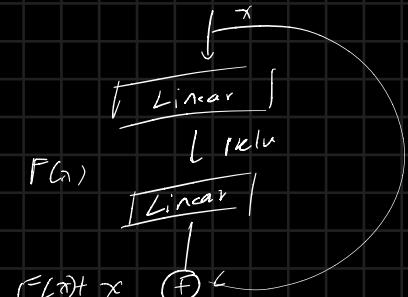
x = torch.save({ "model\_state\_dict": model.state\_dict(), "optimizer\_state\_dict": optimizer.state\_dict() }, "model.pth")

## Skip / Residual Connections

Take output, skip layers, then add to a hidden layer.

$\Rightarrow$  Reduces impact of the vanishing gradient

Input  $\oplus$  hidden must be of same size to add



Creating modules in modules

e.g. create separate skip/residual) blocks

nn. BatchNorm 2d

nn. Sequential

nn. Dropout

Class Conv Block (nn. Module):

def \_\_init\_\_(self, channels):

super(Conv Block, self).\_\_init\_\_()

self.conv1 = nn.Conv2d(channels, channels // 2, kernel\_size=3, stride=1, padding=1)

self.bn1 = nn.BatchNorm2d(channels // 2)

self.conv2 = nn.Conv2d(channels // 2, channels, kernel\_size=3, stride=1, padding=1)

def forward(self, x):

x = F.relu(x)

x = F.relu(self.bn1(self.conv1(x)))

x = self.conv2(x)

return x

Class Res Block (nn. Module):

def \_\_init\_\_(self, channels):

self.conv1 = nn.Conv2d(...)

self.bn1 = nn.BatchNorm2d(...)

self.conv2 = nn.Conv2d(...)

def forward(self, x):

x0 = x

x = F.relu(x)

x = F.relu(self.bn1(self.conv1(x)))

x = self.conv2(x)

return x + x0

Focus down Block (nn.Module):

```
def __init__(self, in_channels, out_channels):
    super(...)
```

self.conv1 = nn.Conv2d(in\_channels, out\_channels, ...)

self.conv2 = nn.Conv2d(out\_channels, out\_channels, ...)

self.bn1 = nn.Conv2d(in\_channels, out\_channels, ...)

self.conv3 = nn.Conv2d(in\_channels, out\_channels, ...)

```
def forward(self, x)
```

$x = F.relu(x)$

$x_0 = self.conv3(x)$

$x = F.relu(self.bn1(self.conv1(x)))$

$x = self.conv2(x)$

return  $x + x_0$

Class Skip Block (nn.Module)

```
def __init__(...)
```

super(...

conv1 (channels, channels//2, kernel\_size=3 ...)

bn1

conv2

```
def forward(self, x)
```

$x = P.relu(x)$

$x_1 = self.conv1(x)$

$x_2 = F.relu(self.bn1(self.conv1(x)))$

$x_3 = self.conv2(x_2)$

return torch.cat([x, x3], 1) # Across depth

Class Deep CNN (nn.Module):

```
def __init__(self, in_channels, num_blocks=2, ch_width=32, layer_type=RcS Block):  
    super().__init__()
```

```
    conv1  
    bn1  
    conv2  
    bn2  
    conv3  
    maxpool
```

self.layers = self.create\_blocks(num\_blocks, layer\_type, channels=ch\_width)

else: nn.Identity()

```
dropout  
linear(...)
```

blocks.append(blocks\_type(channels))

return nn.Sequential(\*blocks)

↑  
unpacks list

def forward(self, x):

x = F.relu(self.bn1(self.conv1(x)))

:

x = F.relu(self.layers(x))

:

:

nn.Sequential([layer1, layer2, ...])

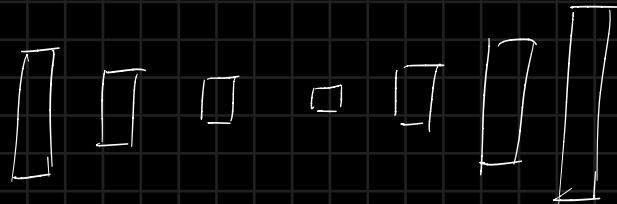
Augmentations:

- T. Auto Augment (?)
- 7. Random Rotation
- 7. Random horizontal flip

Learning rate Scheduler

model.set\_lr\_scheduler(optim.lr\_scheduler.CosineAnnealingLR(optimizer, T\_max=num\_epochs, eta\_min=0))

Auto encoders



Down Block

conv1 => (in, out, 3, 2, 1)

s = 2

p = 1

Up Block

conv1 (in, out, 3, 1, 1)

ConvTranspose2d (in, out, 3, 2, 1, 1)

ConvTranspose2d (...) is the opposite of Conv2d (...)

ConvTranspose2d (in, out, kernel, stride, padding, output\_padding)

Auto encoders are useful for denoising

KL divergence :  $p(x) = N(\mu_p, \sigma_p)$   
 $q(x) = N(\mu_q, \sigma_q)$

$$KL(p||q) = \ln\left(\frac{\sigma_q}{\sigma_p}\right) + \frac{\sigma_p^2 + (\mu_p - \mu_q)^2}{2\sigma_q^2} - \frac{1}{2}$$

Variational auto encoders : represents complex data distributions in a latent space

Latent variables  $z$

Image  $\rightarrow z$ ,  $z \rightarrow$  output  
Encoder      Decoder

VAE : make dist learned by encoder close

to Gaussian, using Kullback-Leibler divergence

Minimize Negative Evidence Lower Bound

$$ELBO : Q(x) = E_{q(z|x)} [\log p(x|z)] - KL(q(z|x) || p(z))$$

KL between 2 normal dists :

$$KL(p||q) \\ KL(p||q) = \ln\left(\frac{\sigma_q}{\sigma_p}\right) + \frac{\sigma_p^2 + (\mu_p - \mu_q)^2}{2\sigma_q^2} - \frac{1}{2}$$

Set  $q(z)$  to be Gaussian

$$KL(q(z|x) || p(z)) = \ln\left(\frac{1}{\sigma}\right) + \frac{\sigma^2 + \mu^2}{2} - \frac{1}{2}$$

KL allows for regularization, and a meaningful encoder structure