

APS360 Fundamentals of AI

Lisa Zhang

Lecture 4; Jan 17, 2019

Agenda

- ▶ Training Terminology
- ▶ Training Curve
- ▶ Hyperparameters
- ▶ Validation Set
- ▶ Assignment 2

Agenda

- ▶ Training Terminology
- ▶ Training Curve
- ▶ Hyperparameters
- ▶ Validation Set
- ▶ Assignment 2

Reminder: Assignment 1 is due Sunday 9pm

Neural Network Training

Training

Terms from last time:

- ▶ Loss function $L(\textit{actual}, \textit{predicted})$
- ▶ Optimizer
- ▶ Training set
- ▶ Test set

Loss (Loss Function)

A **loss function** $L(actual, predicted)$ computes how “bad” a set of predictions was, compared to the ground truth.

- ▶ Large loss = the network’s prediction differs from the ground truth
- ▶ Small loss = the network’s prediction matches the ground truth

~~Q: What are the inputs to the loss function?~~

what does the loss function depend on?

Optimizer

An optimizer determines, based on the value of the **loss function**, how each parameter should change.

The optimizer solves the **credit assignment problem**: how do we assign credit (blame) to the parameters when the network performs poorly?

Optimize Step

We take **one step** towards solving the optimization problem:

$$\min_{weights} L(actual, predicted, weights)$$

How do we do this?

Optimize Step

We take **one step** towards solving the optimization problem:

$$\min_{weights} L(actual, predicted, weights)$$

How do we do this?

Using an optimizer like **gradient descent**.

Optimizer: Gradient Descent

All neural network optimizers you see in this course will be based on **gradient descent**.

We use the derivative of the loss function at a training example, and take a step towards its negative gradient.

You don't need to know how optimizers work for this course.

From learning to optimization

Defining a loss function turned a **learning problem** into an **optimization problem**.

- ▶ Recurrent theme in Machine Learning

Caveats



Custard Smingleigh

@Smingleigh

Follow



I hooked a neural network up to my Roomba. I wanted it to learn to navigate without bumping into things, so I set up a reward scheme to encourage speed and discourage hitting the bumper sensors.

It learnt to drive backwards, because there are no bumpers on the back.

classify edibl
age of the d
n't actually

grow r

ing program, g
side it output a
n by the evalu
output of sort
the empty s

lid mo
to run

Jim Stormdancer @mogwai_poet

Someone compiled a list of instances of AI doing what creators specify, not what they mean:

docs.google.com/spreadsheets/u...

Show this thread

1:18 AM - 8 Nov 2018

5,280 Retweets 13,116 Likes



Train, and Test Set

- ▶ **Training Set:** Used to tune parameters
- ▶ **Test Set:** Used to measure network accuracy

Training and Test Splits

For standard data sets, there are standard train/test splits:

```
mnist_train = datasets.MNIST('data', train=True)
mnist_test = datasets.MNIST('data', train=False)
```

Why?

Code from last week

```
for (image, label) in mnist_train[:1000]:  
    # actual ground truth: is the digit less than 3?  
    actual = (label < 3).reshape([1,1]) \  
                .type(torch.FloatTensor)  
  
    # prediction  
    out = pigeon(img_to_tensor(image))  
    # update the parameters based on the loss  
    loss = criterion(out, actual) # compute loss  
    loss.backward()             # compute param updates  
    optimizer.step()           # make param updates  
    optimizer.zero_grad()      # clean up
```

Summary of Code

1. use our network to make the predictions for **one image**
2. compute the loss for that **one image**
3. take a “step” to optimize the loss of the **one image**

Batching

1. use our network to make the predictions for n **images**
2. compute the *average* loss for those n **image**
3. take a “step” to optimize the *average* loss of those n **image**

Averaging Loss

- ▶ Average loss across multiple training inputs is less “noisy”
- ▶ Less likely to provide “bad information” because of a single “bad” input

Batching Code

```
train_loader = torch.utils.data.DataLoader(  
    mnist_train,  
    batch_size=64)  
for n, (imgs, labels) in enumerate(train_loader):  
    if n >= 10: break  
    actual = (label < 3).reshape([1,1]) \  
        .type(torch.FloatTensor)  
    out = pigeon(img_to_tensor(image))  
    loss = criterion(out, actual)  
    loss.backward()  
    optimizer.step()  
    optimizer.zero_grad()
```

Batching Code

```
train_loader = torch.utils.data.DataLoader(  
    mnist_train,  
    batch_size=64)  
for n, (imgs, labels) in enumerate(train_loader):  
    if n >= 10: break  
    actual = (label < 3).reshape([1,1]) \  
        .type(torch.FloatTensor)  
    out = pigeon(img_to_tensor(image))  
    loss = criterion(out, actual)  
    loss.backward()  
    optimizer.step()  
    optimizer.zero_grad()
```

The inside of the loop looks exactly the same!

Batch Size

The **batch size** is the number of training examples used per optimization “step”.

Each optimization “step” is known as an **iteration**.

The parameters are updated once in an iteration.

Q: What happens if the batch size is too small? Too large?

Ineffective Batch Size

- ▶ **Too small:**
 - ▶ We optimize a (possibly very) different function L at each iteration
 - ▶ Noisy
- ▶ **Too large:**
 - ▶ Expensive
 - ▶ Average loss might not change very much as batch size grows

Epoch

An **epoch** is a measure of the number of times all training data are used once to update the parameters.

Example:

- ▶ There are 1000 images we use for training
- ▶ If `batch_size = 10` then 100 iterations = 1 epoch

Optimizer Settings

- ▶ The optimizer settings can also affect the speed of neural network training.

```
optimizer = optim.SGD(pigeon.parameters()  
                      lr=0.005,  
                      momentum=0.9)
```


Learning Rate

The **learning rate** determines the size of the “step” that an optimizer takes during each *iteration*.

Larger step size = make a bigger change in the parameters in each iteration.

Q: What happens if the learning rate is small? Large?

Learning Rate Size

- ▶ **Too small:**
 - ▶ Parameters don't change very much in each iteration
 - ▶ Takes a long time to train the network
- ▶ **Too large:**
 - ▶ “Noisy”
 - ▶ Average loss might not change very much as batch size grows
 - ▶ Very large can be detrimental to neural network training

Appropriate Learning Rate

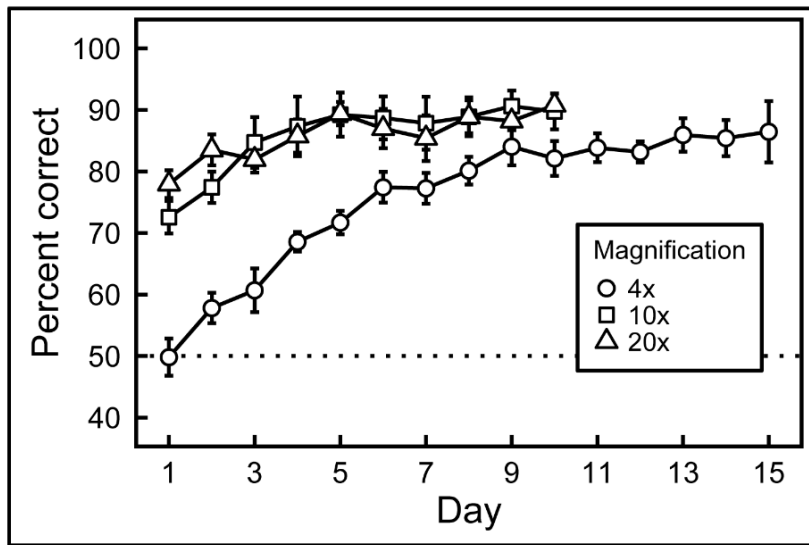
Depends on:

- ▶ The learning problem
- ▶ The optimizer
- ▶ The batch size
 - ▶ Smaller learning rate for larger batch size
 - ▶ Larger learning rate for smaller batch size
- ▶ The stage of training
 - ▶ *Reduce* learning rate as training progresses

Tracking Training

- ▶ How do we know when to stop training?
- ▶ Is training going well?
- ▶ Do we have a good batch size?
- ▶ Do we have a good learning rate?

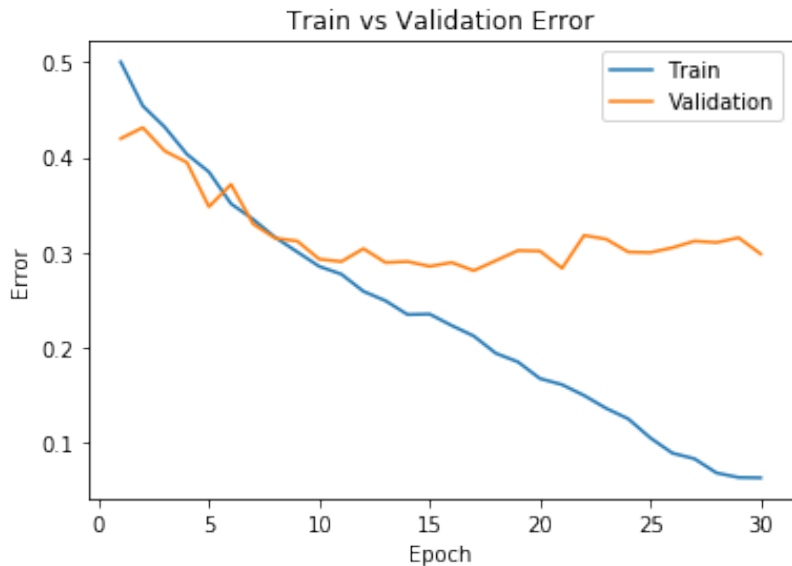
Training Curve for Biological Pigeon



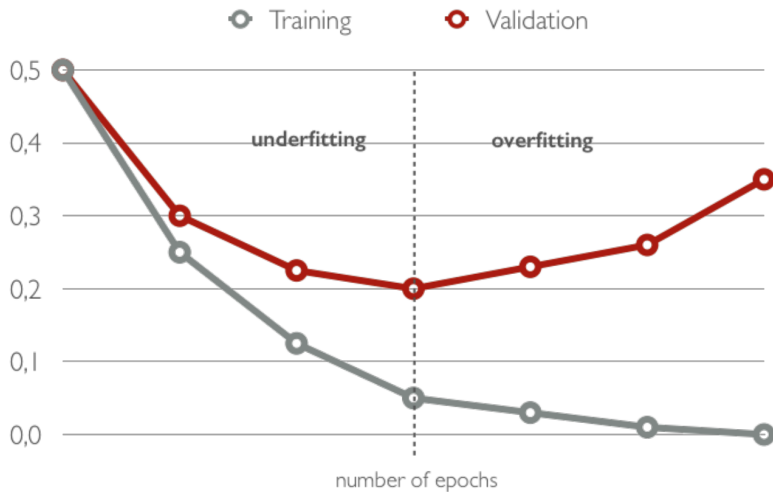
Training Curve

- ▶ **x-axis:** epochs or iterations
- ▶ **y-axis:** loss, error, or accuracy

Typical Training Curve



Assessing the Fit



Hyperparameters

- ▶ Size of network
 - ▶ Number of layers
 - ▶ Number of neurons in each layer
- ▶ Choice of Activation Function
- ▶ Learning Rate
- ▶ Batch Size

Q: How do we tune hyperparameters?

Assignment 2

- ▶ Distinguishing cats and dogs
- ▶ You have pretty much everything you need to begin assignment 2!