Deep Learning Overview

Robert Platt Northeastern University



Problem we want to solve

<u>Scenario</u>:

- A pattern exists
- We don't know what it is, but we have a bunch of examples of the pattern

<u>Machine Learning problem</u>: find a rule for making predictions from the data

Examples:

Image classification: Given a dataset of images and corresponding class labels, learn a rule for predicting the image class given the image

Spam detection: Given a dataset of email text and corresponding spam/ham labels, learn a rule for predicting whether spam/ham given the text of a new email

Problem we want to solve

<u>Scenario</u>:

- A pattern exists
- We don't know what it is, but we have a bunch of examples of the pattern

<u>Machine Learning problem</u>: find a rule for making predictions from the data

<u>Classification vs regression</u>:

- if a labels are discrete, then we have a *classification* problem
- if the labels are real-valued, then we have a *regression* problem

Problem we want to solve

Input: x

Label: y

Data:
$$\mathbb{D} = \{(x_1, y_1), \dots, (x_n, y_n)\}$$

Given ${\mathbb D}$, find a rule for predicting x given y

A single unit neural network



A single unit neural network



Different activation functions:

- id $f(z) = \frac{1}{1 + e^{-z}}$ $f(z) = tanh(z) = \frac{e^z e^{-z}}{e^z + e^{-z}}$ – sigmoid
- tanh

- rectified linear unit $f(z) = \max(0, z)$ (ReLU)



A single unit neural network



One-layer neural network has a simple interpretation: linear classification.

$$y = f(w^T x + b)$$

X_1 == symmetry X_2 == avg intensity Y == class label (binary)



Think-pair-share

$$y = f(w^T x + b)$$

X_1 == symmetry X_2 == avg intensity Y == class label (binary)



What do *w* and *b* correspond to in this picture?

Given a dataset: $D = \{(x^1, y^1), (x^2, y^2), \dots, (x^n, y^n)\}$

Define loss function: $L(x^i, y^i; w, b) = \frac{1}{2}(y^i - f(w^T x^i + b))^2$

Given a dataset: $D = \{(x^1, y^1), (x^2, y^2), \dots, (x^n, y^n)\}$

Define loss function: $L(x^i, y^i; w, b) = \frac{1}{2}(y^i - f(w^T x^i + b))^2$

Loss function tells us how well the network classified x^i

Given a dataset: $D = \{(x^1, y^1), (x^2, y^2), \dots, (x^n, y^n)\}$

Define loss function: $L(x^i, y^i; w, b) = \frac{1}{2}(y^i - f(w^T x^i + b))^2$

Loss function tells us how well the network classified x^i

Method of training: adjust *w*, *b* so as to minimize the net loss over the dataset

i.e.: adjust *w*, *b* so as to minimize:

$$\sum_{(x^i, y^i) \in D} L(x^i, y^i; w, b)$$

The closer to zero, the better the classification

Method of training: adjust *w*, *b* so as to minimize the net loss over the dataset

i.e.: adjust *w*, *b* so as to minimize:

$$\sum_{(x^i,y^i)\in D} L(x^i,y^i;w,b)$$



Method of training: adjust *w*, *b* so as to minimize the net loss over the dataset

i.e.: adjust *w*, *b* so as to minimize:

$$\sum_{(x^i,y^i)\in D} L(x^i,y^i;w,b)$$

How?

Gradient Descent

Gradient descent

Suppose someone gives you an unknown function F(x)

- you want to find a minimum for ${\sf F}$
- but, you do not have an analytical description of F(x)

Use gradient descent!

– all you need is the ability to evaluate F(x) and its gradient at any point x

- 1. pick x_0 at random
- 2. $x_1 = x_0 \alpha \nabla_x F(x_0)$ 3. $x_2 = x_1 - \alpha \nabla_x F(x_1)$ 4. $x_3 = x_2 - \alpha \nabla_x F(x_2)$



5. ...

Gradient descent



1. pick x_0 at random

2. $x_1 = x_0 - \alpha \nabla_x F(x_0)$ 3. $x_2 = x_1 - \alpha \nabla_x F(x_1)$ 4. $x_3 = x_2 - \alpha \nabla_x F(x_2)$



5. ...

Think-pair-share



Training a one-unit neural network



Going deeper: a one layer network



Each hidden node is connected to every input

Going deeper: a one layer network



Going deeper: a one layer network



This is called "forward propogation"

Think-pair-share



Vector of hidden layer activations

Write an expression for *y* in terms of *x*, *f*, and the weights

Hint: use matrix multiplication

Can create networks of arbitrary depth...



- Forward propagation works the same for any depth network.
- Whereas a single output node corresponds to linear classification, adding hidden nodes makes classification non-linear

How do we train multi-layer networks?

Almost the same as in the single-node case...

Do gradient descent on dataset:

1. repeat

2.
$$w \leftarrow w - \alpha \frac{1}{n} \sum_{\substack{(x^i, y^i) \in D}} \nabla_w L(x^i, y^i; w, b)$$

3. $b \leftarrow b - \alpha \frac{1}{n} \sum_{\substack{(x^i, y^i) \in D}} \nabla_b L(x^i, y^i; w, b)$

4. until converged

Now, we're doing gradient descent on all weights/biases in the network – not just a single layer

- this is called *backpropagation*

Stochastic gradient descent: mini-batches

A batch is typically between 32 and 128 samples

1. repeat

2. randomly sample a mini-batch: $B \subset D$

3.
$$w \leftarrow w - \alpha \frac{1}{n} \sum_{\substack{(x^i, y^i) \in D}} \nabla_w L(x^i, y^i; w, b)$$

4. $b \leftarrow b - \alpha \frac{1}{n} \sum_{\substack{(x^i, y^i) \in D}} \nabla_b L(x^i, y^i; w, b)$

5. until converged

Training in mini-batches helps b/c:

- don't have to load the entire dataset into memory
- training is still relatively stable
- random sampling of batches helps avoid local minima

Deep multi-layer perceptron networks

- general purpose
- involve huge numbers of weights

We want:

- special purpose network for image and NLP data
- fewer parameters
- fewer local minima

Answer: convolutional layers!







Because of the way weights are tied together

- reduces number of parameters (dramatically)
- encodes a prior on structure of data

In practice, convolutional layers are essential to computer vision...

Two dimensional example:



Why do you think they call this "convolution"?

Think-pair-share



What would the convolved feature map be for this kernel?





MNIST dataset: images of 10,000 handwritten digits

Objective: classify each image as the corresponding digit

LeNet:



Load dataset, create train/test splits

Define the neural network structure:

```
layers = [
    imageInputLayer([28 28 1])
    convolution2dLayer(3,16, 'Padding',1)
                                                    Input
    batchNormalizationLayer
    reluLayer
    maxPooling2dLayer(2, 'Stride',2)
    convolution2dLayer(3,32, 'Padding',1)
                                                   Conv1
    batchNormalizationLayer
    reluLayer
    maxPooling2dLayer(2, 'Stride',2)
    convolution2dLayer(3,64, 'Padding',1)
                                                   Conv2
    batchNormalizationLayer
    reluLayer
                                                    FC1
    fullyConnectedLayer(50)
    reluLayer
                                                    FC2
    fullyConnectedLayer(10)
    softmaxLayer
    classificationLayer];
options = trainingOptions('sgdm',...
'MaxEpochs',3, ...
'ValidationData', valDigitData,...
'ValidationFrequency',30,...
'Verbose', true, ...
'ExecutionEnvironment', 'gpu',...
'Plots', 'training-progress');
```

Train network, classify test set, measure accuracy

- notice we test on a different set (a holdout set) than we trained on

net = trainNetwork(trainDigitData,layers,options);

```
predictedLabels = classify(net,valDigitData);
valLabels = valDigitData.Labels;
```

accuracy = sum(predictedLabels == valLabels)/numel(valLabels);



Using the GPU makes a huge differece...

Deep learning packages

You don't need to use Matlab (obviously)

Tensorflow is probably the most popular platform

Caffe and Theano are also big



Caffe

theano



ImageNet dataset: millions of images of objects

Objective: classify each image as the corresponding object (1k categories in ILSVRC)



AlexNet has 8 layers: five conv followed by three fully connected

1	'data'	Image Input	227x227x3 images with 'zerocenter' normalization
2	'conv1'	Convolution	96 11x11x3 convolutions with stride [4 4] and padding [0 0]
3	'relu1'	ReLU	ReLU
4	'norm1'	Cross Channel Normalization	cross channel normalization with 5 channels per element
5	'pool1'	Max Pooling	3x3 max pooling with stride [2 2] and padding [0 0]
6	'conv2'	Convolution	256 5x5x48 convolutions with stride [1 1] and padding [2 2]
7	'relu2'	ReLU	ReLU
8	'norm2'	Cross Channel Normalization	cross channel normalization with 5 channels per element
9	'pool2'	Max Pooling	3x3 max pooling with stride [2 2] and padding [0 0]
10	'conv3'	Convolution	384 3x3x256 convolutions with stride [1 1] and padding [1 1]
11	'relu3'	ReLU	ReLU
12	'conv4'	Convolution	384 3x3x192 convolutions with stride [1 1] and padding [1 1]
13	'relu4'	ReLU	ReLU
14	'conv5'	Convolution	256 3x3x192 convolutions with stride [1 1] and padding [1 1]
15	'relu5'	ReLU	ReLU
16	'pool5'	Max Pooling	3x3 max pooling with stride [2 2] and padding [0 0]
17	'fc6'	Fully Connected	4096 fully connected layer
18	'relu6'	ReLU	ReLU
19	'drop6'	Dropout	50% dropout
20	'fc7'	Fully Connected	4096 fully connected layer
21	'relu7'	ReLU	ReLU
22	'drop7'	Dropout	50% dropout
23	'fc8'	Fully Connected	1000 fully connected layer
24	'prob'	Softmax	softmax
25	'output'	Classification Output	crossentropyex with 'tench', 'goldfish', and 998 other classes

AlexNet has 8 layers: five conv followed by three fully connected

AlexNet won the 2012 ILSVRC challenge – sparked the deep learning craze

