Photogrammetry & Robotics Lab

Intro to Neural Networks Part 2: Learning

Cyrill Stachniss

The slides have been created by Cyrill Stachniss.

5 Minute Preparation for Today



https://www.ipb.uni-bonn.de/5min/

5 Minute Preparation for Today



https://www.ipb.uni-bonn.de/5min/

In Part 1, We Discussed

- What are neurons and neural networks
- Activations, weights, biases
- Multi-layer perceptron (MLP)
- MLP for simple image classification



Part 2 Learning the Parameters

How to Make the Network Compute What We Want?

- Neural network is a recipe for performing a set of computations
- Structure and parameters are the design choices
- How to set them?



Network Parameters

Given a network structure, weights and biases tell the network what to do



params = weights & biases

$$\boldsymbol{W}^{(1)} \; \boldsymbol{b}^{(1)} \; \cdots \; \boldsymbol{W}^{(k)} \; \boldsymbol{b}^{(k)}$$

What Means "Learning"?

- NN are functions
- The weights and biases determine what this function computes



- Learning = determining parameters so that the network does what we want
- Parameters are estimated by providing labeled examples (training data)

Our Handwritten Digit Network



simple combined patterns raw pixels patterns patterns to digits

[Image courtesy: Nielsen]

Parameters Encode Features



Many Parameters

Such networks have **many** parameters!



Training Through Labeled Data



Training Through Labeled Data



Exploiting Training Examples



Loss Function

 We define a loss (= cost) function over all weights and biases of the network

$L(W, \boldsymbol{b}) \mapsto \mathbb{R}$

- Computed using training data
- Input to *L* are the network parameters $\boldsymbol{\theta} = \boldsymbol{W}, \boldsymbol{b}$
- Output is the error of the network with these parameters on the training data

Loss Function $L_i(\boldsymbol{\theta}) \mapsto \mathbb{R}$

 $\hat{\boldsymbol{y}}_i = f(\boldsymbol{\theta}, \boldsymbol{x}_i) \quad \boldsymbol{y}_i$ \boldsymbol{x}_i 0.2Input Laver Hidden Laver 1 Hidden Laver 2 Output Laver 784 128 64 10 0 0 (relu) (relu) (softmax) 0.30 0.20 0.20 0.91 0 0 network with label data network point parameters θ output Compare output layer to the true label $L_i(\boldsymbol{\theta}) = ||\text{output}_i - \text{label}_i||^2$

Loss Over All Examples

We need to evaluate the performance of the network over all examples

loss = all examples, avg. squared(NN(input) - label)



The Parameters We Want

• Parameter θ^* that minimize the sum of avg. squared losses over all examples

$$oldsymbol{ heta}^* = rg\min_{oldsymbol{ heta}} L\left(oldsymbol{ heta}
ight) = rg\min_{oldsymbol{ heta}} \sum_i ||f(oldsymbol{ heta}, oldsymbol{x}_i) - oldsymbol{y}_i||^2$$

- The squared loss is only one possible loss, several other options available
- Goal: Find the parameter vector θ*
 for the labeled training set {(x_i, y_i)}^I_{i=1}
 given the loss L

Let's Start...

- Initialize parameters randomly
- See how well it performs (bad!)
- How to improve the parameters so that the loss decreases?

$$L(\boldsymbol{\theta}) \mapsto \mathbb{R}$$

high 1 dim
dimensional
(109.386 dim)

That's complex!



Loss Minimization using Gradient Descent

- Our problem looks like a non-linear least squares problem
- We have a lot of parameters, which makes using GN computationally tricky
- Gradient descent is a better way to perform the minimization



How to move?



How to move? Exploit the gradient



23





 $\nabla L = \frac{\partial L}{\partial \theta}$ tells us in which direction to go



Step by step: $\theta^{(j+1)} = \theta^{(j)} - \lambda \nabla L|_{\theta^{(j)}}$

- First derivative of loss: $\nabla L = \frac{\partial L}{\partial \theta}$
- Learning rate (small value): $\lambda = 0.01$

Gradient descent works by

- $\theta^{(0)} = \operatorname{rand}$
- while (!converged) $\theta^{(j+1)} = \theta^{(j)} - \lambda |_{\theta^{(j)}}$



Step by step: $\theta^{(j+1)} = \theta^{(j)} - \lambda \nabla L|_{\theta^{(j)}}$

- We can do the same in 2D
- Gradients are direction vectors



- We can do the same in 2D
- Gradients are direction vectors



 \bigcirc

- We can do the same in 2D
- Gradients are direction vectors



- We can do the same in 2D
- Gradients are direction vectors



Meaning of the Gradient Vector

- Some dimensions are more important than others to reduce the loss
- Gradient indicates which change leads to the fastest reduction of the loss



changes in θ_2 are more relevant than changes in θ_1 to reduce loss

 θ_2



changes in θ_1 and θ_2 have the same relevance to reduce loss

Gradient Descent in Higher Dimensional Spaces

- Same situation as before, but the gradient vector has more dimensions
- Update rule

$$\boldsymbol{\theta}^{(j+1)} = \boldsymbol{\theta}^{(j)} - \lambda \nabla L|_{\boldsymbol{\theta}^{(j)}}$$

$$\uparrow \qquad \uparrow \qquad \uparrow$$

In our classification example, all these vectors have 109.386 dimensions!

Keep in Mind...



We need to adjust the parameters to minimize the total loss!

total loss is the averaged loss of all examples

$$L\left(\boldsymbol{\theta}\right) = \frac{1}{I} \sum_{i} L_{i}\left(\boldsymbol{\theta}\right)$$

Gradient Over All Examples

Loss function sums over all examples

$$L\left(\boldsymbol{\theta}\right) = \frac{1}{I} \sum_{i} L_{i}\left(\boldsymbol{\theta}\right)$$

This means for the gradient

$$\nabla L = \frac{1}{I} \sum_{i} \nabla L_{i}$$

We need to sum over all training examples and compute all gradients whenever we perform a single GD step!
Two Challenges

1. How to optimize the process if we have a **lot of training examples**?

2. How to compute the gradients for **complex and nested functions**?

We need to perform these operations often, so we need to be able to execute them efficiently!

1: Handling Large Training Sets

1st trick: Compute a gradient only on a small, sampled subset of examples



1: Handling Large Training Sets

1st trick: Compute a gradient only on a small, sampled subset of examples



1: Stochastic Gradient Descent

- 1st trick: Compute a gradient only on a small, sampled subset of examples
- We sample a mini-batch in each step of gradient descent
- Use only mini-batch B to compute

$$\nabla L = \frac{1}{|B|} \sum_{i \in B} \nabla L_i$$

This approximates the real gradient

1: Stochastic Gradient Descent

Approximate down-hill steps but much faster to compute



2: Computing the Gradient

- 2nd trick: Compute ∇L_i step by step
- Neuron activations are chains of activation functions and matrix-vector multiplications
- Many connections between the neurons
- Computing this 109.386 dimensional gradient can be tricky...

backpropagation algorithm

2: Backpropagation

 The idea is to break down the gradient computation into smaller steps



Key ingredients of backpropagation:

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial x}$$
chain rule

local variables along paths through the NN

Backpropagation

Computational Graph

- Directed graph
- Nodes contain mathematical operations
- Edges encode input/output values
- Example:



Function and Corresponding Computational Graph

$$f(x, y, z) = (x + y) z^2$$



Evaluating the Function

$$f(x, y, z) = (x + y) z^2 \qquad f(2, -3, 2) = -4$$



Add Local Variables

 $f(x, y, z) = (x + y) z^2$



Computing the Gradient

- With the forward pass, we evaluate the function at a given point
- Next: compute the gradient of the function at the given point
- We can do this by traversing the graph backwards
- At each note, we compute the local derivative of the function w.r.t. the local inputs

Gradient At Point [2, -3, 2]





























Backpropagation (Backprop)

- Approach is called backpropagation and computes the gradient of any function expressed as such a graph
- Combines chain rule and local variables for the graph nodes
- Recursively traverses the graph
- Can be used for computing both, numerical and analytical gradient computation

Example for One Neuron

$$a = \sigma \left(\sum w_n a_n + b \right) \quad \sigma(x) = \frac{1}{1 + \exp(-x)}$$



Example for One Neuron with Two Incoming Activations



Forward Pass



Forward Pass










Backward Pass











Backward Pass



Backpropagation for a Single Neuron with Sigmoid Activ.

- The example illustrated that we can compute the gradient for a neuron
- We can also model the sigmoid using a single node ("sigmoid gate")



Backpropagation from Neuron to Neuron

- Recursive application from neuron to neuron through the layers
- For multiple outgoing edges, we need to compute the sum of gradients



 Backpropagation also generalizes directly to multivariate function

Neural Networks As Computation Graphs

- We can use BP for computing the gradient of the loss function $L(\theta)$
- We add a loss layer



Quadratic Loss for Backpropagation

training data



parameters

Backpropagation (Backprop)

- BP allows us to compute gradients of nested and complex functions
- Combines chain rule and local variables
- Recursive traversal of the network
- Forward pass computes the linearization point for the gradient
- Backward pass computes the gradient

Learning a Neural Network

Repeat until convergence

- 1. Sample mini-batch from training data
- Run backprop to compute gradient for SGD using mini-batch

$$\nabla L|_{\boldsymbol{\theta}^{(j)}}$$

3. Execute SGD step to find better parameters reducing the loss

$$\boldsymbol{\theta}^{(j+1)} = \boldsymbol{\theta}^{(j)} - \lambda \nabla L|_{\boldsymbol{\theta}^{(j)}}$$

Return parameter vector

In image-related learning tasks, CNNs play an important role



[Image courtesy: van Veen] 85

In image-related learning tasks, CNNs play an important role



In image-related learning tasks, CNNs play an important role



In image-related learning tasks, CNNs play an important role



[Image courtesy: van Veen] 88

What Is "Deep Learning"?

"Learning neural networks with many hidden layers"

Definition: #hidden layers > 2

Very deep networks today have up to 150 hidden layers (still growing...)

Summary – Part 2

- Leaning multi-layer perceptrons
- Parameters are the weights and biases
- Learning = estimate weights & biases
- Minimization of a loss (cost) function
- Gradient descent for parameter optimization
- Backpropagation to compute gradients
- End-to-end: no manual features
- CNN for image processing

Literature & Resources

- Online Book by Michael Nielsen, Chapter 1: http://neuralnetworksanddeeplearning.com/chap1.html
- Nielsen, Chapter 1, Python3 code: https://github.com/MichalDanielDobrzanski/DeepLearningPython
- MNIST database:
- http://yann.lecun.com/exdb/mnist/
- Grant Sanderson, Neural Networks https://www.3blue1brown.com/
- Online book by Deisenroth, Faisal, Ong: Mathematics for Machine Learning https://mml-book.github.io/
- Alpaydin, Introduction to Machine Learning
- Standford AI Lectures by Li et al.

Slide Information

- The slides have been created by Cyrill Stachniss as part of the photogrammetry and robotics courses.
- I tried to acknowledge all people from whom I used images or videos. In case I made a mistake or missed someone, please let me know.
- Huge thank you to Grant Sanderson (3blue1brown) for his great educational videos that influenced this lecture.
- Thanks to Michael Nielsen for his free online book & code as well as to Fei-Fei Li et al. for the Stanford AI lectures.
- If you are a university lecturer, feel free to use the course material. If you adapt the course material, please make sure that you keep the acknowledgements to others and please acknowledge me as well. To satisfy my own curiosity, please send me email notice if you use my slides.

Cyrill Stachniss, cyrill.stachniss@igg.uni-bonn.de