ML4 B&I: Introduction to Machine Learning Lecture 5- Neural Networks

Murat A. Erdogdu

Vector Institute, Fall 2022

Outline

1 Limits of Linear Classification

2 Introducing Neural Networks

3 Backpropagation

Recall: Visualizing NOT



• Data is linearly separable if a linear decision rule can perfectly separate the training examples.

Some datasets are not linearly separable, e.g. XOR.



Classifying XOR Using Feature Maps

Sometimes, we can overcome this limitation using feature maps, e.g., for **XOR**.

		x_1	x_2	$\psi_1(\mathbf{x})$	$\psi_2(\mathbf{x})$	$\psi_3(\mathbf{x})$	t
$oldsymbol{\psi}(\mathbf{x}) = egin{pmatrix} x_1 \ x_2 \ x_1 x_2 \end{pmatrix}$	_	0	0	0	0	0	0
		0	1	0	1	0	1
)	1	0	1	0	0	1
	,	1	1	1	1	1	0

- This is linearly separable. (Try it!)
- Designing feature maps can be hard. Can we learn them?





3 Backpropagation

Neurons in the Brain

Neurons receive input signals and accumulate voltage. After some threshold, they will fire spiking responses.



[Pic credit: www.moleculardevices.com]

A Simpler Neuron

For neural nets, we use a much simpler model for neuron, or **unit**:



A Simpler Neuron

For neural nets, we use a much simpler model for neuron, or **unit**:



- Same as logistic regression: $y = \sigma(\mathbf{w}^{\top}\mathbf{x} + b)$
- By throwing together lots of these simple neuron-like processing units, we can do some powerful computations!

A Feed-Forward Neural Network



an output unit

Multilayer Perceptrons

- A multi-layer network consists of fully connected layers.
- In a fully connected layer, all input units are connected to all output units.

Multilayer Perceptrons

- A multi-layer network consists of fully connected layers.
- In a fully connected layer, all input units are connected to all output units.
- The outputs are a function of the input units:

$$\mathbf{y} = f(\mathbf{x}) = \phi\left(\mathbf{W}\mathbf{x} + \mathbf{b}\right)$$

 ϕ is applied component-wise.



Some Activation Functions



More Activation Functions



Computation in Each Layer

Each layer computes a function.

$$\begin{aligned} \mathbf{h}^{(1)} &= f^{(1)}(\mathbf{x}) = \phi(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) \\ \mathbf{h}^{(2)} &= f^{(2)}(\mathbf{h}^{(1)}) = \phi(\mathbf{W}^{(2)}\mathbf{h}^{(1)} + \mathbf{b}^{(2)}) \end{aligned}$$

$$\begin{array}{c|c} T & \bigcirc & \bigcirc \\ \hline f^{(L)} \\ \hline \\ & \\ f^{(3)} \\ \hline \\ 2) & \bigcirc & \bigcirc \\ \end{array}$$

J

)



•

 $\mathbf{y} = f^{(L)}(\mathbf{h}^{(L-1)})$

Computation in Each Layer

Each layer computes a function.

$$\begin{aligned} \mathbf{h}^{(1)} &= f^{(1)}(\mathbf{x}) = \phi(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) \\ \mathbf{h}^{(2)} &= f^{(2)}(\mathbf{h}^{(1)}) = \phi(\mathbf{W}^{(2)}\mathbf{h}^{(1)} + \mathbf{b}^{(2)}) \end{aligned}$$

$$\mathbf{y} = f^{(L)}(\mathbf{h}^{(L-1)})$$

•

- If task is regression: choose $\mathbf{y} = f^{(L)}(\mathbf{h}^{(L-1)}) = (\mathbf{w}^{(L)})^{\top} \mathbf{h}^{(L-1)} + b^{(L)}$
- If task is binary classification: choose $\mathbf{y} = f^{(L)}(\mathbf{h}^{(L-1)}) = \sigma((\mathbf{w}^{(L)})^{\top}\mathbf{h}^{(L-1)} + b^{(L)})$



A Composition of Functions

The network computes a composition of functions.

$$\mathbf{y} = f^{(L)} \circ \cdots \circ f^{(1)}(\mathbf{x}).$$

Modularity: We can implement each layer's computations as a black box.



Feature Learning

Neural nets can be viewed as a way of learning features:



The goal:



Feature Learning



- Suppose we're trying to classify images of handwritten digits.
- Each image is represented as a vector of $28 \times 28 = 784$ pixel values.
- Each hidden unit in the first layer acts as a **feature detector**.
- We can visualize **w** by reshaping it into an image. Below is an example that responds to a diagonal stroke.



Features for Classifying Handwritten Digits

Features learned by the first hidden layer of a handwritten digit classifier:



Unlike hard-coded feature maps (e.g., in polynomial regression), features learned by neural networks adapt to patterns in the data.

- Consider a linear layer: the activation function was the identity. The layer just computes an affine transformation of the input.
- Any sequence of linear layers is equivalent to a single linear layer.

$$\mathbf{y} = \underbrace{\mathbf{W}^{(3)}\mathbf{W}^{(2)}\mathbf{W}^{(1)}}_{\triangleq \mathbf{W}'} \mathbf{x}$$

• Deep linear networks can only represent linear functions — no more expressive than linear regression.

Expressive Power of Non-linear Networks

- Multi-layer feed-forward neural networks with non-linear activation functions
- Universal Function Approximators: They can approximate any function arbitrarily well.
- True for various activation functions (e.g. thresholds, logistic, ReLU, etc.)

Designing a Network to Classify XOR

Assume a hard threshold activation function.



Designing a Network to Classify XOR

 h_1 is computed as $x_1 \vee x_2$

$$h_1 = \mathbb{I}[x_1 + x_2 - 0.5 > 0]$$

 h_2 is computed as $x_1 \wedge x_2$

$$h_2 = \mathbb{I}[x_1 + x_2 - 1.5 > 0]$$

y computes

$$y = \mathbb{I}[h_1 - h_2 - 0.5 > 0]$$

= x_1 XOR x_2



Expressivity of the Logistic Activation Function

- What about the logistic activation function?
- Approximate a hard threshold by scaling up w and b.



Expressivity of the Logistic Activation Function

- What about the logistic activation function?
- Approximate a hard threshold by scaling up w and b.



• Logistic units are differentiable, so we can learn weights with gradient descent.

What is Expressivity Good For?

- May need a very large network to represent a function.
- Non-trivial to learn the weights that represent a function.
- If you can learn any function, over-fitting is potentially a serious concern!

For the polynomial feature mappings, expressivity increases with the degree M, eventually allowing multiple perfect fits to the training data. This motivated L^2 regularization.



• Do neural networks over-fit and how can we regularize them?

ML4 B&I-Lec5

Regularization and Over-fitting for Neural Networks

- The topic of over-fitting (when & how it happens, how to regularize, etc.) for neural networks is not well-understood, even by researchers!
 - ▶ In principle, you can always apply L^2 regularization.
- A common approach is early stopping, or stopping training early, because over-fitting typically increases as training progresses.



- **1** Limits of Linear Classification
- **2** Introducing Neural Networks
- 3 Backpropagation

Learning Weights in a Neural Network

- Goal is to learn weights in a multi-layer neural network using gradient descent.
- Weight space for a multi-layer neural net: one set of weights for each unit in every layer of the network
- Define a loss \mathcal{L} and compute the gradient of the cost $d\mathcal{J}/d\mathbf{w}$, the average loss over all the training examples.
- Let's look at how we can calculate $d\mathcal{L}/d\mathbf{w}$.

Example: Two-Layer Neural Network

Figure: Two-Layer Neural Network

Intro ML (Vector)

ML4 B&I-Lec5

Computations for Two-Layer Neural Network

A neural network computes a composition of functions.

$$z_1^{(1)} = w_{01}^{(1)} \cdot 1 + w_{11}^{(1)} \cdot x_1 + w_{21}^{(1)} \cdot x_2$$

$$h_1 = \sigma(z_1)$$

$$z_1^{(2)} = w_{01}^{(2)} \cdot 1 + w_{11}^{(2)} \cdot h_1 + w_{21}^{(2)} \cdot h_2$$

$$y_1 = z_1$$

$$z_2^{(1)} =$$

$$h_2 =$$

$$z_2^{(2)} =$$

$$y_2 =$$

$$L = \frac{1}{2} \left((y_1 - t_1)^2 + (y_2 - t_2)^2 \right)$$

Simplified Example: Logistic Least Squares

- The nodes represent the inputs and computed quantities.
- The edges represent which nodes are computed directly as a function of which other nodes.

Let z = f(y) and y = g(x) be uni-variate functions. Then z = f(g(x)).

$$\frac{\mathrm{d}z}{\mathrm{d}x} = \frac{\mathrm{d}z}{\mathrm{d}y} \ \frac{\mathrm{d}y}{\mathrm{d}x}$$

Logistic Least Squares: Gradient for w

Computing the loss:

$$z = wx + b$$

$$y = \sigma(z)$$

$$\mathcal{L} = \frac{1}{2}(y - t)^{2}$$

Computing the gradient for w:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial w} &= \frac{\partial \mathcal{L}}{\partial y} \frac{\partial y}{\partial w} \\ &= \frac{\partial \mathcal{L}}{\partial y} \frac{\partial y}{\partial z} \frac{\partial z}{\partial w} \\ &= (y-t) \sigma'(z) x \\ &= (\sigma(wx+b)-t)\sigma'(wx+b)x \end{aligned}$$

Intro ML (Vector)

Logistic Least Squares: Gradient for b

Computing the loss:

$$z = wx + b$$

$$y = \sigma(z)$$

$$\mathcal{L} = \frac{1}{2}(y - t)^{2}$$

Computing the gradient for b:

$$\frac{\partial \mathcal{L}}{\partial b} =$$

Logistic Least Squares: Gradient for b

Computing the loss:

$$z = wx + b$$

$$y = \sigma(z)$$

$$\mathcal{L} = \frac{1}{2}(y - t)^{2}$$

Computing the gradient for b:

$$\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial y} \frac{\partial y}{\partial b}$$
$$= \frac{\partial \mathcal{L}}{\partial y} \frac{\partial y}{\partial z} \frac{\partial z}{\partial b}$$
$$= (y - t) \sigma'(z) 1$$
$$= (\sigma(wx + b) - t)\sigma'(wx + b)1$$

Intro ML (Vector)

Comparing Gradient Computations for w and b

Computing the loss:

$$z = wx + b$$

$$y = \sigma(z)$$

$$\mathcal{L} = \frac{1}{2}(y - t)^{2}$$

Computing the gradient for w: Computing the gradient for b:

Intro ML (Vector)

Structured Way of Computing Gradients

Computing the loss:

$$z = wx + b$$

$$y = \sigma(z)$$

$$\mathcal{L} = \frac{1}{2}(y - t)^{2}$$

Computing the gradients:

$$\frac{\partial \mathcal{L}}{\partial y} = (y - t)$$
$$\frac{\partial \mathcal{L}}{\partial z} = \frac{\partial \mathcal{L}}{\partial y} \ \sigma'(z)$$

 $\frac{\partial \mathcal{L}}{\partial w} = \frac{\mathrm{d}\mathcal{L}}{\mathrm{d}z}\frac{\mathrm{d}z}{\mathrm{d}w} = \frac{\mathrm{d}\mathcal{L}}{\mathrm{d}z}x \qquad \qquad \frac{\partial \mathcal{L}}{\partial b} = \frac{\mathrm{d}\mathcal{L}}{\mathrm{d}z}\frac{\mathrm{d}z}{\mathrm{d}b} = \frac{\mathrm{d}\mathcal{L}}{\mathrm{d}z}1$ Intro ML (Vector) ML4 B&I-Lec5

- Let \overline{y} denote the derivative $d\mathcal{L}/dy$, called the **error signal**.
- Error signals are just values our program is computing (rather than a mathematical operation).

Computing the loss:

Computing the derivatives:

$$z = wx + b$$

$$y = \sigma(z)$$

$$\mathcal{L} = \frac{1}{2}(y - t)^{2}$$

$$\overline{y} = (y - t)$$
$$\overline{z} = \overline{y} \, \sigma'(z)$$
$$\overline{w} = \overline{z} \, x \qquad \overline{b} = \overline{z}$$

Computation Graph has a Fan-Out > 1

L_2 -Regularized Regression

Computation Graph has a Fan-Out > 1

Softmax Regression

Multi-variate Chain Rule

Suppose we have functions f(x, y), x(t), and y(t).

$$\frac{\mathrm{d}}{\mathrm{d}t}f(x(t), y(t)) = \frac{\partial f}{\partial x}\frac{\mathrm{d}x}{\mathrm{d}t} + \frac{\partial f}{\partial y}\frac{\mathrm{d}y}{\mathrm{d}t}$$

Example:

$$\begin{aligned} f(x,y) &= y + e^{xy} & \frac{\mathrm{d}f}{\mathrm{d}t} = \frac{\partial f}{\partial x}\frac{\mathrm{d}x}{\mathrm{d}t} + \frac{\partial f}{\partial y}\frac{\mathrm{d}y}{\mathrm{d}t} \\ x(t) &= \cos t & \\ y(t) &= t^2 & = (ye^{xy}) \cdot (-\sin t) + (1 + xe^{xy}) \cdot 2t \end{aligned}$$

Multi-variate Chain Rule

In the context of back-propagation:

In our notation:

$$\overline{t} = \overline{x} \, \frac{\mathrm{d}x}{\mathrm{d}t} + \overline{y} \, \frac{\mathrm{d}y}{\mathrm{d}t}$$

ML4 B&I-Lec5

Backpropagation for Regularized Logistic Least Squares

Forward pass:

$$z = wx + b$$

$$y = \sigma(z)$$

$$\mathcal{L} = \frac{1}{2}(y - t)^{2}$$

$$\mathcal{R} = \frac{1}{2}w^{2}$$

$$\mathcal{L}_{reg} = \mathcal{L} + \lambda \mathcal{R}$$

Backpropagation for Regularized Logistic Least Squares

Backpropagation for Regularized Logistic Least Squares

 $\overline{z} = \overline{y} \frac{\mathrm{d}y}{\mathrm{d}z}$ $\overline{\mathcal{R}} = \overline{\mathcal{L}_{\mathrm{reg}}} \frac{\mathrm{d}\mathcal{L}_{\mathrm{reg}}}{\mathrm{d}\mathcal{R}}$ $=\overline{\mathcal{L}_{\mathrm{reg}}}\lambda$ $= \overline{y} \sigma'(z)$ $\overline{w} = \overline{z} \frac{\partial z}{\partial w} + \overline{\mathcal{R}} \frac{\mathrm{d}\mathcal{R}}{\mathrm{d}w}$ $\overline{\mathcal{L}} = \overline{\mathcal{L}_{\mathrm{reg}}} \, \frac{\mathrm{d}\mathcal{L}_{\mathrm{reg}}}{\mathrm{d}\ell}$ $=\overline{z}x+\overline{\mathcal{R}}w$ $=\overline{\mathcal{L}_{reg}}$ $\overline{b} = \overline{z} \frac{\partial z}{\partial b}$ $\overline{y} = \overline{\mathcal{L}} \, \frac{\mathrm{d}\mathcal{L}}{\mathrm{d}y}$ $=\overline{7}$ $=\overline{\mathcal{L}}(y-t)$

Computational Cost

• Computational cost of forward pass: one add-multiply operation per weight

$$z_i = \sum_j w_{ij}^{(1)} x_j + b_i^{(1)}$$

• Computational cost of backward pass: two add-multiply operations per weight

$$\overline{w_{ki}^{(2)}} = \overline{y_k} h_i$$
$$\overline{h_i} = \sum_k \overline{y_k} w_{ki}^{(2)}$$

• One backward pass is as expensive as two forward passes.

- The algorithm for efficiently computing gradients in neural nets.
- Gradient descent with gradients computed via backprop is used to train the overwhelming majority of neural nets today.
- Even optimization algorithms much fancier than gradient descent (e.g. second-order methods) use backprop to compute the gradients.

- Autodifferentiation performs backprop in a completely mechanical and automatic way.
- Many autodiff libraries: PyTorch, Tensorflow, Jax, etc.
- Although autodiff automates the backward pass for you, it's still important to know how things work under the hood.
- In the assignment, you will use an autodiff framework to build complex neural networks.

- Introduced Neural Networks
- Discuss their expressive power.
 - Can approximate any function.
- Introduced backpropagation.
 - ▶ In the appendix, we also work out the updates for a two-layer neural network.

Appendix: Full Backpropagation Algorithm:

Let v_1, \ldots, v_N be an ordering of the computation graph where parents come before children.

 v_N denotes the variable for which we're trying to compute gradients.

• forward pass:

For i = 1, ..., N, Compute v_i as a function of Parents (v_i) .

• backward pass:

For
$$i = N - 1, \dots, 1$$
,
 $\bar{v}_i = \sum_{j \in \text{Children}(v_i)} \bar{v}_j \frac{\partial v_j}{\partial v_i}$

Appendix: Backpropagation for Two-Layer Neural Network

Forward pass:

$$z_{i} = \sum_{j} w_{ij}^{(1)} x_{j} + b_{i}^{(1)}$$
$$h_{i} = \sigma(z_{i})$$
$$y_{k} = \sum_{i} w_{ki}^{(2)} h_{i} + b_{k}^{(2)}$$
$$\mathcal{L} = \frac{1}{2} \sum_{k} (y_{k} - t_{k})^{2}$$

Backward pass:

$$\overline{\mathcal{L}} = 1$$

$$\overline{y_k} = \overline{\mathcal{L}} (y_k - t_k)$$

$$\overline{w_{ki}^{(2)}} = \overline{y_k} h_i$$

$$\overline{b_k^{(2)}} = \overline{y_k}$$

$$\overline{h_i} = \sum_k \overline{y_k} w_{ki}^{(2)}$$

$$\overline{z_i} = \overline{h_i} \sigma'(z_i)$$

$$\overline{w_{ij}^{(1)}} = \overline{z_i} x_j$$

$$\overline{b_i^{(1)}} = \overline{z_i}$$

Intro ML (Vector)

ML4 B&I-Lec5

Appendix: Backpropagation for Two-Layer Neural Network

In vectorized form:

Forward pass:

$$\begin{aligned} \mathbf{z} &= \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)} \\ \mathbf{h} &= \sigma(\mathbf{z}) \\ \mathbf{y} &= \mathbf{W}^{(2)}\mathbf{h} + \mathbf{b}^{(2)} \\ \mathcal{L} &= \frac{1}{2}\|\mathbf{t} - \mathbf{y}\|^2 \end{aligned}$$

Backward pass:

 $\overline{f} = 1$ $\overline{\mathbf{v}} = \overline{\mathcal{L}} \left(\mathbf{v} - \mathbf{t} \right)$ $\overline{\mathbf{W}^{(2)}} = \overline{\mathbf{y}}\mathbf{h}^{\top}$ $\overline{\mathbf{b}^{(2)}} = \overline{\mathbf{v}}$ $\overline{\mathbf{h}} = \mathbf{W}^{(2)\top} \overline{\mathbf{y}}$ $\overline{\mathbf{z}} = \overline{\mathbf{h}} \circ \sigma'(\mathbf{z})$ $\overline{\mathbf{W}^{(1)}} = \overline{\mathbf{z}}\mathbf{x}^{\top}$ $\overline{\mathbf{b}^{(1)}} = \overline{\mathbf{z}}$

Intro ML (Vector)

ML4 B&I-Lec5