

CSCI-567: Machine Learning

Prof. Victor Adamchik

U of Southern California

July 9, 2020

Your model is only as good as your data.

July 9, 2020 1 / 46

Outline

- 1 Problem Solving
- 2 Perceptron
- 3 From Perceptron to Neural Network
- 4 Backpropagation

July 9, 2020 3 / 46

Outline

- 1 Problem Solving
- 2 Perceptron
- 3 From Perceptron to Neural Network
- 4 Backpropagation

July 9, 2020 2 / 46

Problem 3

Show that one-versus-all can be seen as a special case of error-correcting-output-code (ECOC). Specifically, write down the code matrix M for ECOC for a problem with C labels so that executing ECOC is the same as doing one-versus-all. (Note: the entry of M should be either -1 or $+1$.)

Problem 4

Assume we have a training set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$, the probability of seeing outcome y is given by

$$P(y|\mathbf{x}_n) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(y - \mathbf{w}^T \mathbf{x}_n)^2}{2\sigma^2}\right)$$

Find the maximum likelihood estimations for \mathbf{w} and σ

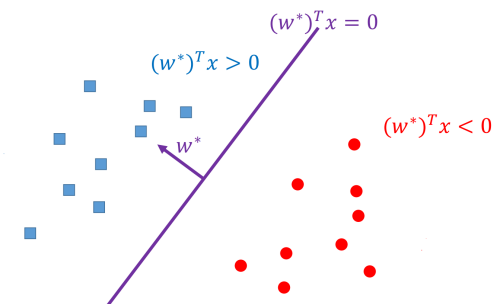
Outline

- 1 Problem Solving
- 2 **Perceptron**
- 3 From Perceptron to Neural Network
- 4 Backpropagation

Solution

The Perceptron

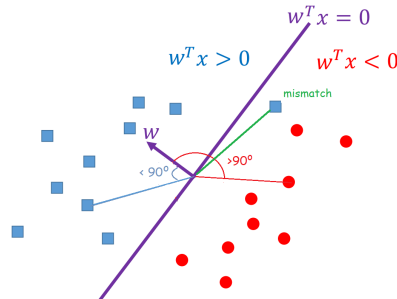
The Perceptron (introduced by Rosenblatt in 1957) is a linear model for classification. Its model is a hyperplane that partitions space into two regions.



Intuition

The product of two vectors is defined by

$$\mathbf{w}^T \mathbf{x} = \|\mathbf{w}\| \|\mathbf{x}\| \cos(\alpha)$$



For a point to be classified as a blue square we want $\mathbf{w}^T \mathbf{x} > 0$, and the angle $\alpha < 90^\circ$ between \mathbf{x} and \mathbf{w} . In order to decrease the alpha value, we update weights by $\mathbf{w}' = \mathbf{w} + \mathbf{x}$

Why does it make sense?

If the current weight \mathbf{w} makes a mistake

$$y_n \mathbf{w}^T \mathbf{x}_n < 0$$

then after the update $\mathbf{w}' = \mathbf{w} + y_n \mathbf{x}_n$ we have

$$y_n \mathbf{w}'^T \mathbf{x}_n = y_n \mathbf{w}^T \mathbf{x}_n + y_n^2 \mathbf{x}_n^T \mathbf{x}_n = y_n \mathbf{w}^T \mathbf{x}_n + \|\mathbf{x}_n\|^2 \geq y_n \mathbf{w}^T \mathbf{x}_n$$

Thus it is more likely to get it right after the update.

The Perceptron Algorithm

Repeat:

- Pick a data point \mathbf{x}_n uniformly at random
- If $\text{sgn}(\mathbf{w}^T \mathbf{x}_n) \neq y_n$

$$\mathbf{w} \leftarrow \mathbf{w} + y_n \mathbf{x}_n$$

Note:

- The algorithm is online and error driven.
- If the prediction is correct, it does nothing.
- \mathbf{w} is always a *linear combination* of the training examples.

Applying GD to perceptron loss

Mathematically, find the minimizer of

$$F(\mathbf{w}) = \sum_{n=1}^N \ell_{\text{perceptron}}(y_n \mathbf{w}^T \mathbf{x}_n) = \sum_{n=1}^N \max\{0, -y_n \mathbf{w}^T \mathbf{x}_n\}$$

Gradient is

$$\nabla F(\mathbf{w}) = \sum_{n=1}^N -\mathbb{I}[y_n \mathbf{w}^T \mathbf{x}_n \leq 0] y_n \mathbf{x}_n$$

(only misclassified examples contribute to the gradient)

GD update

$$\mathbf{w} \leftarrow \mathbf{w} + \lambda \sum_{n=1}^N \mathbb{I}[y_n \mathbf{w}^T \mathbf{x}_n \leq 0] y_n \mathbf{x}_n$$

Any theory?

- If training set is linearly separable, Perceptron *converges in a finite number of steps*
- How long does it take to converge?
- By "how long", what we really mean is "how many updates".
- One way to make this definition is through the notion of margin.
- The margin is the distance between the hyperplane and the nearest point.

Perceptron Convergence Theorem

Suppose the perceptron algorithm is run on a linearly separable data set \mathcal{D} with margin $\gamma \geq 0$. Assume that $\|x\| = 1$. Then the algorithm will converge after at most $\frac{1}{\gamma^2}$ updates.

Exercise 1

The following table shows a binary classification training set and the number of times each point is misclassified during a run of the perceptron algorithm. Which of the following is the final output of the algorithm? Assume w initialized with zeros.

x	y	Times misclassified
$(-3, 2)$	$+1$	5
$(-1, 1)$	-1	5
$(5, 2)$	$+1$	3
$(2, 2)$	-1	4
$(1, -2)$	$+1$	3

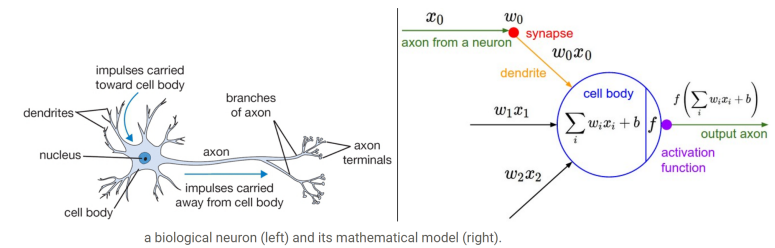
Outline

- 1 Problem Solving
- 2 Perceptron
- 3 From Perceptron to Neural Network
- 4 Backpropagation

Biological motivation

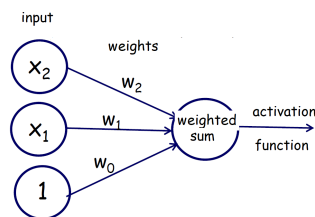
Perceptron is a rough model for how individual neurons in the brain work. The basic computational unit of the brain is a **neuron**. Approximately 86 billion neurons can be found in the human brain. Neurons are connected with each other via **synapses**. Each neuron receives input signals from its **dendrites** and produces output signals along its **axon**.

Biological motivation



In the computational model of a neuron, the signals that travel along the axons (e.g. x_0) interact with the dendrites of the other neuron (e.g. w_0x_0). In the basic model, the dendrites carry the signal to the cell body where they all get summed. If the final sum is above a certain threshold, the neuron can **fire**. We model the firing rate of the neuron by an activation function.

Perceptron as a linear classifier

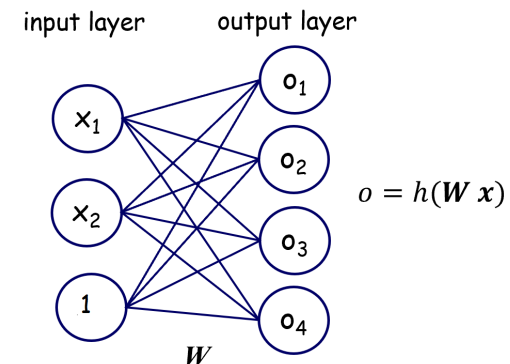


The perceptron works in three simple steps:

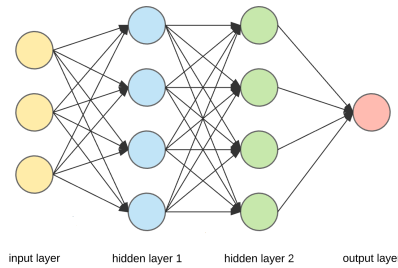
- All the inputs x are multiplied with their weights w .
- All the multiplied values are summed to weighted sum.
- An activation function is applied to that weighted sum.

An activation function must be a nonlinear function!

More output nodes



More layers



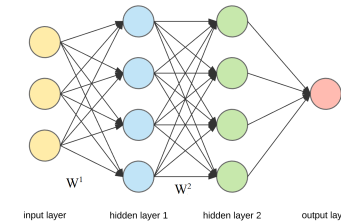
Neural Network:

- each node is called a **neuron**
- each hidden layer has an **activation function**
- #layers refers to #hidden_layers (plus 1 or 2 for input/output layers)
- **deep** neural nets can have many layers and *millions* of parameters
- this is a **feedforward, fully connected** neural net, called multilayer perceptron (MLP).

Math formulation

An L-layer neural net can be written as

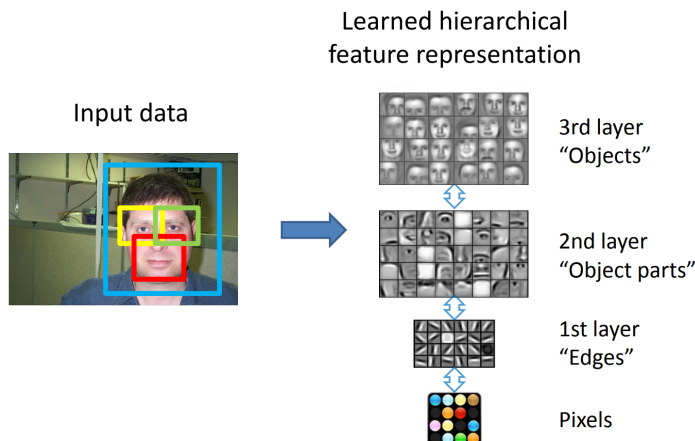
$$f(x) = h_L(W^L h_{L-1}(W^{L-1} \dots h_1(W^1 x)))$$



To learn this model we need to learn weights W^1, W^2, \dots, W^L . They are learnt by Gradient Descent using *the backpropagation algorithm*.

Example from Honglak Lee (NIPS 2010)

Face recognition:



How powerful are neural nets?

Universal approximation theorem (Cybenko, 89; Hornik, 91):

A feedforward neural net with a single hidden layer and finite number of neurons can approximate any continuous functions under mild assumptions on the activation (sigmoidal) function.

It might need a huge number of neurons though, and *depth helps!* With more layers, you may need less neurons.

Designing network architecture is important and very complicated

- for feedforward network, need to decide number of hidden layers, number of neurons at each layer, activation functions, etc.

There is no theory for that, but some rules of thumb in practice.

Deep Feedforward Networks: Historical Notes

In the 1940s, approximation techniques were used to motivate machine learning models such as the perceptron.

Inability to learn the XOR function, led to a backlash against the neural network.

Following the success of backpropagation (1986), NN gained popularity.

The core ideas behind modern feedforward networks have not changed substantially since the 1980s.

The same backpropagation algorithm and the same approaches to gradient descent are still in use.

The modern revival of deep learning started in 2006.

Until 2012, it was widely believed that feedforward networks would not perform well.

Today, it is now known that with the right resources and engineering practices, feedforward networks perform very well.

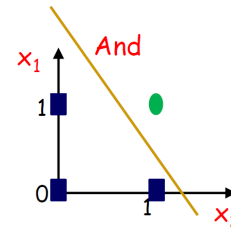
July 9, 2020 22 / 46

Solution

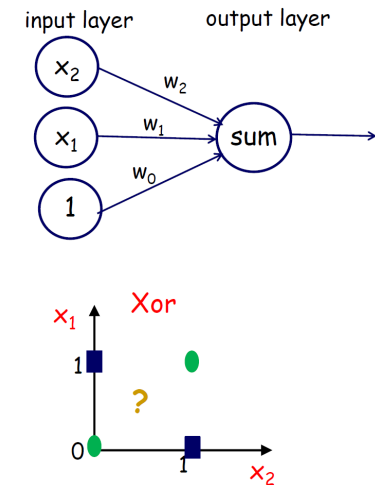
July 9, 2020 24 / 46

Exercise 2

Perceptron can represent some Boolean functions.



Prove why perceptron cannot represent XOR.



July 9, 2020 23 / 46

Outline

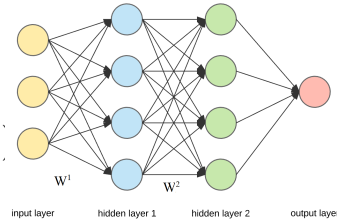
- 1 Problem Solving
- 2 Perceptron
- 3 From Perceptron to Neural Network
- 4 **Backpropagation**
 - Definition
 - Backpropagation
 - Preventing overfitting

July 9, 2020 25 / 46

Math formulation

An L-layer neural net can be written as

$$f(x) = h_L(W^L h_{L-1}(W^{L-1} \dots h_1(W^1 x)))$$



To ease notation, for a given input x , (one training point) define recursively

$$o^0 = x, \quad a^\ell = W^\ell o^{\ell-1}, \quad o^\ell = h_\ell(a^\ell) \quad (\ell = 1, \dots, L)$$

where

- $D_0 = D, D_1, \dots, D_L$ are numbers of neurons at each layer
- $W^\ell \in \mathbb{R}^{D_\ell \times D_{\ell-1}}$ is the weights for layer ℓ
- $a^\ell \in \mathbb{R}^{D_\ell}$ is input to layer ℓ
- $o^\ell \in \mathbb{R}^{D_\ell}$ is output to layer ℓ
- $h_\ell : \mathbb{R}^{D_\ell} \rightarrow \mathbb{R}^{D_\ell}$ is activation functions at layer ℓ

July 9, 2020 26 / 46

Learning the model

No matter how complicated the model is, our goal is the same: minimize

$$\mathcal{E}(W_1, \dots, W_L) = \sum_{n=1}^N \mathcal{E}_n(W_1, \dots, W_L)$$

where

$$\mathcal{E}_n(W_1, \dots, W_L) = \begin{cases} \|f(x_n) - y_n\|_2^2 & \text{for regression} \\ \ln \left(1 + \sum_{k \neq y_n} e^{f(x_n)_k - f(x_n)_{y_n}} \right) & \text{for classification} \end{cases}$$

July 9, 2020 27 / 46

How to optimize such a complicated function?

Same thing: apply **SGD**! even if the model is *nonconvex*.

What is the gradient of this complicated function?

Chain rule:

- for a composite function $f(g(w))$

$$\frac{\partial f}{\partial w} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial w}$$

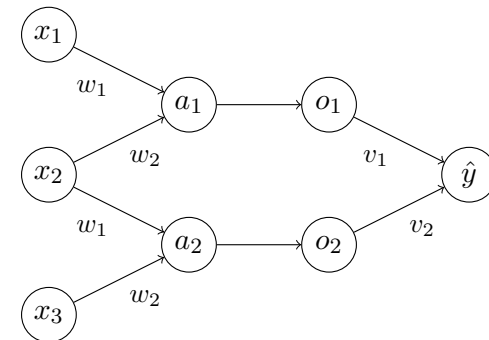
- for a composite function $f(g_1(w), \dots, g_d(w))$

$$\frac{\partial f}{\partial w} = \sum_{i=1}^d \frac{\partial f}{\partial g_i} \frac{\partial g_i}{\partial w}$$

July 9, 2020 28 / 46

Simple example

Consider the following mini neural net, where (x_1, x_2, x_3) is the input, followed by a layer with a filter (w_1, w_2) , a ReLU layer, and another layer with weight (v_1, v_2) .



where $(x, y) \in \mathbb{R}^3 \times \{-1, +1\}$, and \hat{y} is prediction, y is ground truth.

July 9, 2020 29 / 46

Simple example

More concretely, the computation is specified by

$$\begin{aligned}a_1 &= x_1 w_1 + x_2 w_2, & a_2 &= x_2 w_1 + x_3 w_2 \\o_1 &= \max\{0, a_1\}, & o_2 &= \max\{0, a_2\} \\ \hat{y} &= o_1 v_1 + o_2 v_2\end{aligned}$$

We define a loss function by

$$\ell = \ln(1 + \exp(-y\hat{y}))$$

Simple example

Next we compute $\frac{\partial \ell}{\partial w_1}$ and $\frac{\partial \ell}{\partial w_2}$. We will be using $H(a) = \mathbb{I}[a > 0]$.

$$\begin{aligned}\frac{\partial \ell}{\partial w_1} &= \frac{\partial \ell}{\partial a_1} \frac{\partial a_1}{\partial w_1} + \frac{\partial \ell}{\partial a_2} \frac{\partial a_2}{\partial w_1} \\ &= \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_1} \frac{\partial o_1}{\partial a_1} \frac{\partial a_1}{\partial w_1} + \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial o_2} \frac{\partial o_2}{\partial a_2} \frac{\partial a_2}{\partial w_1} \\ &= (\sigma(y\hat{y}) - 1)y(v_1 H(a_1)x_1 + v_2 H(a_2)x_2)\end{aligned}$$

$$\frac{\partial \ell}{\partial w_2} = (\sigma(y\hat{y}) - 1)y(v_1 H(a_1)x_2 + v_2 H(a_2)x_3).$$

Simple example

We start with the last layer and compute $\frac{\partial \ell}{\partial v_1}$ and $\frac{\partial \ell}{\partial v_2}$.

$$\begin{aligned}\frac{\partial \ell}{\partial v_1} &= \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial v_1} \\ &= \frac{-ye^{-y\hat{y}}}{1 + e^{-y\hat{y}}} o_1 = -\sigma(-y\hat{y})y o_1 = (\sigma(y\hat{y}) - 1)y o_1\end{aligned}$$

$$\begin{aligned}\frac{\partial \ell}{\partial v_2} &= \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial v_2} \\ &= \frac{-ye^{-y\hat{y}}}{1 + e^{-y\hat{y}}} o_2 = -\sigma(-y\hat{y})y o_2 = (\sigma(y\hat{y}) - 1)y o_2\end{aligned}$$

Backpropagation algorithm.

Input: A training set $(x_1, y_1), \dots, (x_N, y_N)$, learning rate λ

Initialize: w_1, w_2, v_1, v_2 randomly

While not converged:

- randomly pick an example (x_n, y_n)
- Forward propagation: compute

$$\begin{aligned}a_1 &= x_1 w_1 + x_2 w_2, & a_2 &= x_2 w_1 + x_3 w_2 \\ o_1 &= \max\{0, a_1\}, & o_2 &= \max\{0, a_2\}, & \hat{y} &= o_1 v_1 + o_2 v_2\end{aligned}$$

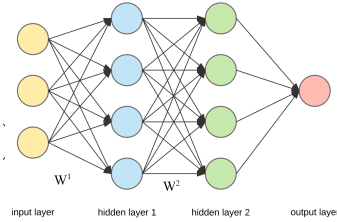
- Backward propagation: update

$$\begin{aligned}w_1 &\leftarrow w_1 - \lambda(\sigma(y\hat{y}) - 1)y(v_1 H(a_1)x_1 + v_2 H(a_2)x_2) \\ w_2 &\leftarrow w_2 - \lambda(\sigma(y\hat{y}) - 1)y(v_1 H(a_1)x_2 + v_2 H(a_2)x_3) \\ v_1 &\leftarrow v_1 - \lambda(\sigma(y\hat{y}) - 1)y o_1 \\ v_2 &\leftarrow v_2 - \lambda(\sigma(y\hat{y}) - 1)y o_2\end{aligned}$$

General Case: math formulation

An L-layer neural net can be written as

$$f(x) = h_L(W^L h_{L-1}(W^{L-1} \dots h_1(W^1 x)))$$



To ease notation, for a given input x , define recursively

$$o^0 = x, \quad a^\ell = W^\ell o^{\ell-1}, \quad o^\ell = h_\ell(a^\ell) \quad (\ell = 1, \dots, L)$$

where

- $D_0 = D, D_1, \dots, D_L$ are numbers of neurons at each layer
- $W^\ell \in \mathbb{R}^{D_\ell \times D_{\ell-1}}$ is the weights for layer ℓ
- $a^\ell \in \mathbb{R}^{D_\ell}$ is input to layer ℓ
- $o^\ell \in \mathbb{R}^{D_\ell}$ is output to layer ℓ
- $h_\ell : \mathbb{R}^{D_\ell} \rightarrow \mathbb{R}^{D_\ell}$ is activation functions at layer ℓ

July 9, 2020 34 / 46

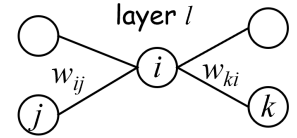
General Case: computing derivatives

An equation for the rate of change of the cost with respect to any weight in the network.

Find the **derivative of \mathcal{E}_n w.r.t. to W**

$$a^l = W^l o^{l-1}$$

$$o^l = h_l(a^l)$$



$$\frac{\partial \mathcal{E}_n}{\partial W^\ell} = \frac{\partial \mathcal{E}_n}{\partial a^\ell} \frac{\partial a^\ell}{\partial W^\ell} = \frac{\partial \mathcal{E}_n}{\partial a^\ell} (o^{\ell-1})^T$$

July 9, 2020 35 / 46

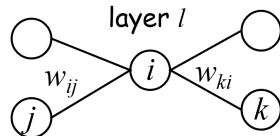
General Case: computing derivatives

Computing the error in terms of the next layer errors.

Find the **derivative of \mathcal{E}_n w.r.t. to a^ℓ**

$$a^l = W^l o^{l-1}$$

$$o^l = h_l(a^l)$$



$$\frac{\partial \mathcal{E}_n}{\partial a_i^\ell} = \frac{\partial \mathcal{E}_n}{\partial o_i^\ell} \frac{\partial o_i^\ell}{\partial a_i^\ell} = \left(\sum_k \frac{\partial \mathcal{E}_n}{\partial a_k^{\ell+1}} \frac{\partial a_k^{\ell+1}}{\partial o_i^\ell} \right) h'_\ell(a_i^\ell) = \left(\sum_k \frac{\partial \mathcal{E}_n}{\partial a_k^{\ell+1}} w_{ki}^{\ell+1} \right) h'_\ell(a_i^\ell)$$

$$\frac{\partial \mathcal{E}_n}{\partial a^\ell} = \left((W^{\ell+1})^T \frac{\partial \mathcal{E}_n}{\partial a^{\ell+1}} \right) \circ h'_\ell(a^\ell)$$

where $x \circ y = (x_1 y_1, \dots, x_n y_n)$ is the Hadamard product.

July 9, 2020 36 / 46

General Case: computing derivatives

Computing the error for the last layer. We know that

$$\mathcal{E}_n(W_1, \dots, W_L) = \begin{cases} \|f(x_n) - y_n\|_2^2 & \text{for regression} \\ \ln \left(1 + \sum_{k \neq y_n} e^{f(x_n)_k - f(x_n)_{y_n}} \right) & \text{for classification} \end{cases}$$

where $f(x) = h_L(a^L)$.

The derivative is :

$$\frac{\partial \mathcal{E}_n}{\partial a^L} = \nabla \mathcal{E}_n \circ h'_L(a^L)$$

July 9, 2020 37 / 46

General Case: backpropagation

Backpropagation is based on the following three equations:

$$\frac{\partial \mathcal{E}_n}{\partial \mathbf{W}^\ell} = \frac{\partial \mathcal{E}_n}{\partial \mathbf{a}^\ell} (\mathbf{o}^{\ell-1})^T$$
$$\frac{\partial \mathcal{E}_n}{\partial \mathbf{a}^\ell} = \begin{cases} ((\mathbf{W}^{\ell+1})^T \frac{\partial \mathcal{E}_n}{\partial \mathbf{a}^{\ell+1}}) \circ \mathbf{h}'_\ell(\mathbf{a}^\ell) & \text{if } \ell < L \\ \nabla \mathcal{E}_n \circ \mathbf{h}'_L(\mathbf{a}^L) & \text{else} \end{cases}$$

where $\mathbf{x} \circ \mathbf{y} = (x_1 y_1, \dots, x_n y_n)$ is the Hadamard product.

Verify dimensions for the above equations.

- $\mathbf{W}^\ell \in \mathbb{R}^{D_\ell \times D_{\ell-1}}$ is the weights for layer ℓ
- $\mathbf{a}^\ell \in \mathbb{R}^{D_\ell}$ is input to layer ℓ
- $\mathbf{o}^\ell \in \mathbb{R}^{D_\ell}$ is output to layer ℓ

Putting everything into SGD

The backpropagation algorithm (1986) computes the gradient of the cost function for a single training example.

Initialize $\mathbf{W}^1, \dots, \mathbf{W}^L$ (all 0 or randomly). Repeat:

- 1 randomly pick one data point $n \in [N]$
- 2 **forward propagation**: for each layer $\ell = 1, \dots, L$
 - compute $\mathbf{a}^\ell = \mathbf{W}^\ell \mathbf{o}^{\ell-1}$ and $\mathbf{o}^\ell = \mathbf{h}_\ell(\mathbf{a}^\ell)$ ($\mathbf{o}^0 = \mathbf{x}_n$)
- 3 **backward propagation**: for each $\ell = L, \dots, 1$

► compute

$$\frac{\partial \mathcal{E}_n}{\partial \mathbf{a}^\ell}$$

► update weights

$$\mathbf{W}^\ell \leftarrow \mathbf{W}^\ell - \lambda \frac{\partial \mathcal{E}_n}{\partial \mathbf{W}^\ell} = \mathbf{W}^\ell - \lambda \frac{\partial \mathcal{E}_n}{\partial \mathbf{a}^\ell} (\mathbf{o}^{\ell-1})^T$$

More tricks to optimize neural nets

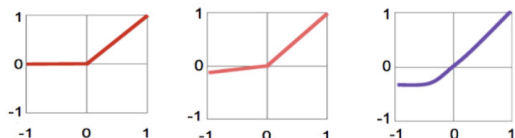
Using ReLU (instead of sigmoid) as an activation function. It solves the vanishing gradient problem.

There are several improved versions of ReLU that have been proposed which provide even better accuracy

One of them is Leaky ReLU: $y = \max(ax, x)$, here a is a small constant, $a = 0.1$.

Another popular activation function is Exponential Linear Unit (ELU) defined by

$$y = \begin{cases} x, & x \geq 0 \\ a(e^x - 1), & x < 0 \end{cases}$$



Overfitting

Overfitting is very likely since the models are too powerful.

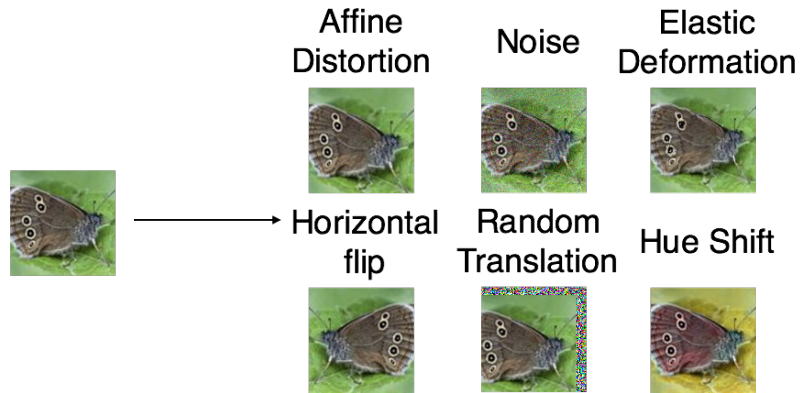
Methods to overcome overfitting:

- data augmentation
- regularization
- dropout
- early stopping
- ...

Data augmentation

Data: the more the better. How do we get more data?

Exploit prior knowledge to add more training data



July 9, 2020 42 / 46

Regularization

L2 regularization: minimize

$$\mathcal{E}'(\mathbf{W}_1, \dots, \mathbf{W}_L) = \mathcal{E}(\mathbf{W}_1, \dots, \mathbf{W}_L) + \eta \sum_{\ell=1}^L \|\mathbf{W}_\ell\|_2^2$$

Simple change to the gradient:

$$\frac{\partial \mathcal{E}'}{\partial w_{ij}} = \frac{\partial \mathcal{E}}{\partial w_{ij}} + 2\eta w_{ij}$$

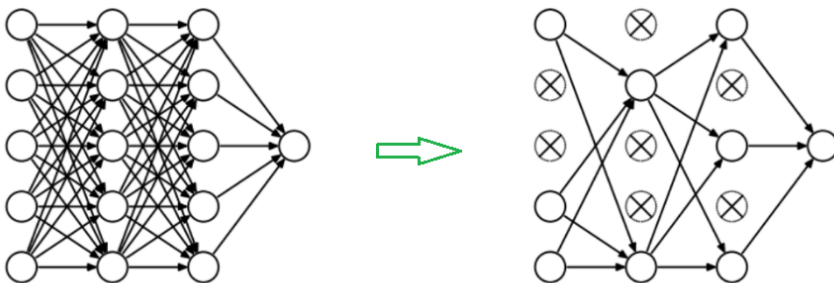
Introduce *weight decaying effect*:

$$\mathbf{W}^\ell \leftarrow \mathbf{W}^\ell - \lambda \frac{\partial \mathcal{E}_n}{\partial \mathbf{W}^\ell} - 2\eta \lambda \mathbf{W}^\ell$$

July 9, 2020 43 / 46

Dropout

Randomly delete neurons during training

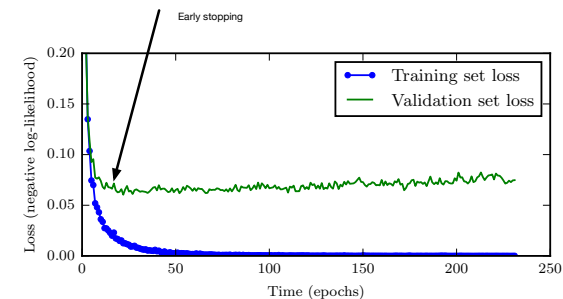


Very effective, makes training faster as well

July 9, 2020 44 / 46

Early stopping

Stop training when the performance on validation set stops improving



July 9, 2020 45 / 46

Conclusions for neural nets

Deep neural networks

- are hugely popular, achieving *best performance* on many problems
- do need *a lot of data* to work well
- take *a lot of time* to train (need GPUs for massive parallel computing)
- take some work to select architecture and hyperparameters
- are still not well understood in theory

Conclusions for backpropagation learning

(due to Yann LeCun)

- gradient descent is unreliable when the network is small
- the solution is to make the network much larger than necessary and regularize it
- Although there are lots of local minima, many of them are equivalent
- It doesn't matter which one you fall into with large nets, backpropagation yields very consistent results
- Many local minima are due to symmetries in the system
- Breaking the symmetry in the architecture solves many problems. This may be why convolutional nets work so well