

Unit 2:

Supervised Learning

How do we adjust a neural network so that it does what we want it to do?
How do we get the network to *learn*?

By the end of this unit, you will be able to...

- Formulate the problem of supervised learning as an optimization problem.
- List and explain some of the most common loss/cost functions.
- Describe a perceptron, and its limitations.
- Use gradient descent to optimize connection weights.
- Derive and implement the error backpropagation algorithm.
- Use labelled data wisely to train models that are generalizable.
- Explain the problem of vanishing/exploding gradients.
- Employ some methods to improve the convergence of our learning method.

Neural Learning

Goal: To formulate the problem of supervised learning as an optimization problem.

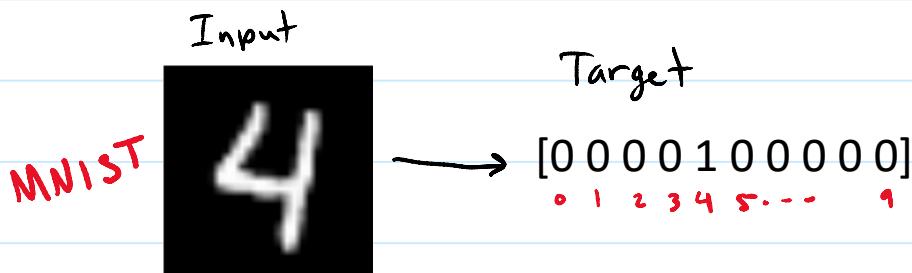
Getting a neural network to do what you want usually means finding a set of connection weights that yield the desired behaviour. That is, neural learning is all about adjusting connection weights.

There are three basic categories of learning problems:

1. Supervised learning
2. Unsupervised learning
3. Reinforcement learning

In **supervised learning**, the desired output is known so we can compute the error and use that error to adjust our network.

Example: Given an image of a digit, identify which digit it is.



In **unsupervised learning**, the output is not known (or not supplied), so cannot be used to generate an error signal. Instead, this form of learning is all about finding efficient representations for the statistical structure in the input.

Example: Given spoken English words, transform them into a more efficient representation such as phonemes, and then syllables.

In **reinforcement learning**, feedback is given, but usually less often, and the error signal is usually less specific.

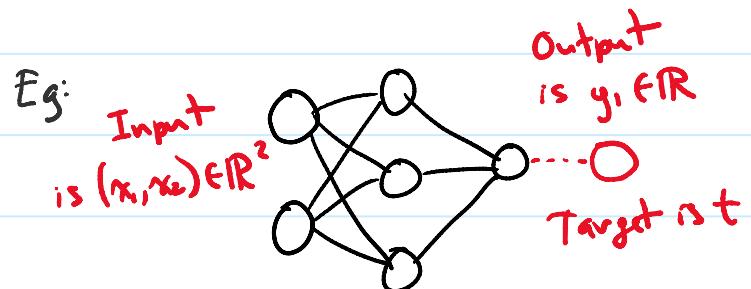
Example: When playing a game of chess, a person knows their play

was good if they win the game. They can try to learn from the moves they made.

In this course, we will mostly focus on supervised learning. But we will look at some examples of unsupervised learning, and possibly some reinforcement learning.

Supervised Learning

Our neural network performs some mapping from an input space to an output space.



We are given training data, with many MANY examples of input/target pairs. This data is (presumably) the result of some consistent mapping process. For example, handwritten digits map to numbers. Or,

A	B	$\text{XOR}(A, B)$
1	1	0
1	0	1
0	1	1
0	0	0

Input $(A, B) \in \{0, 1\}^2$

Output/Target

$y \in \mathbb{R}$ or $y \in [0, 1]$ $t \in \{0, 1\}$

We are given inputs and their corresponding targets, one pair (or a few pairs) at a time. Our task is to alter the connection weights in our network so that our network mimics this mapping.

Our goal is to bring the output as close as possible to the target. But what, exactly, do we mean by "close"? For now, we will use the scalar function $E(y, t)$ as an error function, which returns a smaller value as our outputs are closer to the target.

Two common types of mappings encountered in supervised learning are **regression** and **classification**.

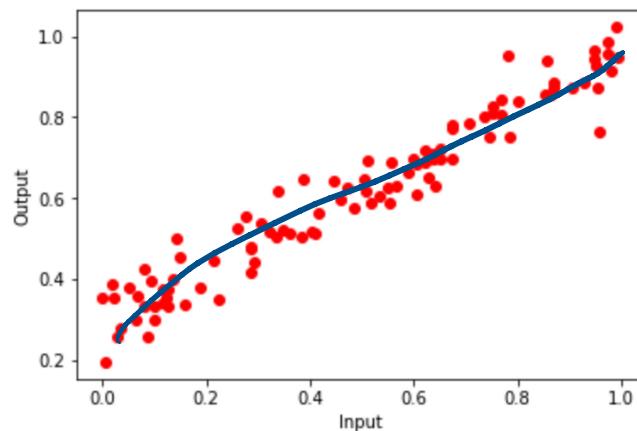
Regression

Output values are a continuous-valued function of the inputs. The outputs can take on a range of values.

Example:

Linear Regression

Outputs fall
range of
values



Classification

Outputs fall into a number of distinct categories.

Example:

MNIST

Inputs

7

[0 0 0 0 0 0 1 0 0]

0

[1 0 0 0 0 0 0 0 0]

6

[0 0 0 0 0 1 0 0 0]

Inputs

5

[0 0 0 0 1 0 0 0 0]

4

[0 0 0 0 1 0 0 0 0]

9

[0 0 0 0 0 0 0 0 1]

CIFAR-10

<u>Inputs</u>	<u>Outputs</u>
	airplane
	automobile
	bird
	cat
	deer
	dog
	frog
	horse
	ship
	truck

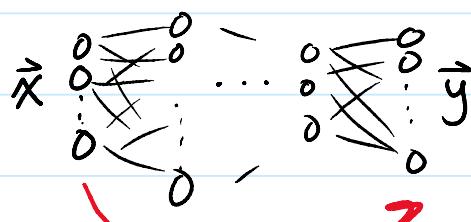
Optimization

Once we have a cost function, our neural-network learning problem can be formulated as an optimization problem.

Let our network be represented by the mapping f , so that

$$\vec{y} = f(\vec{x}; \theta)$$

where θ represents all the weights and biases.



$$\min_{\theta} E \left[E(f(\vec{x}; \theta), \underbrace{f(\vec{x})}_{\text{output of model}}) , \underbrace{\vec{x}}_{\text{target } \vec{x} \in \text{data}} \right]$$

In other words, find the weights and biases that minimize the expected cost between the outputs and the targets.

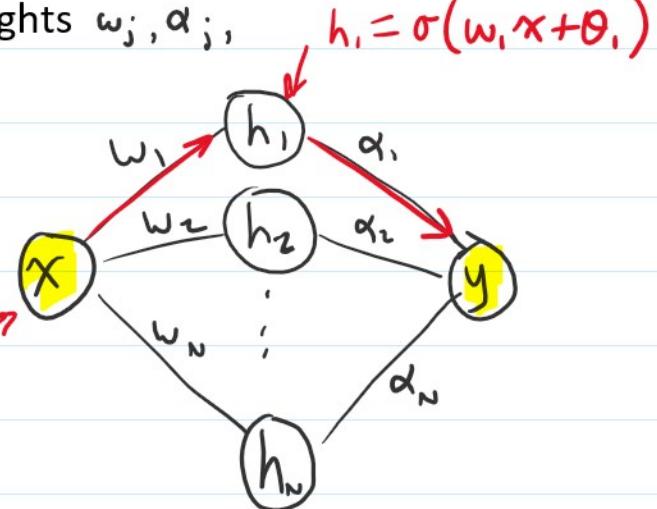
Universal Approximation Theorem

Question: Can we approximate **any** function using a neural network?

Given a function $f(x)$, can we find the weights w_j, α_j , and biases $\theta_j, j=1, \dots, N$ such that

$$f(x) \approx \sum_{j=1}^N \alpha_j \sigma(w_j x + \theta_j)$$

to arbitrary precision?



Universal Approximation Theorem:

Theorem 2. Let σ be any continuous sigmoidal function. Then finite sums of the form

continuous functions

$$G(x) = \sum_{j=1}^N \alpha_j \sigma(w_j x + \theta_j)$$



are dense in $C(I_n)$. In other words, given any $f \in C(I_n)$ and $\epsilon > 0$, there is a sum, $G(x)$, of the above form, for which

$$|G(x) - f(x)| < \epsilon \quad \text{for all } x \in I_n \quad I_n = [0, 1]$$

Cybenko G, "Approximation by Superpositions of a Sigmoidal Function", *Math. Control Signals Systems*, 2:303-314, 1989.

A function σ is "sigmoidal" if $\sigma(x) = \begin{cases} 1 & \text{as } x \rightarrow \infty \\ 0 & \text{as } x \rightarrow -\infty \end{cases}$
e.g. logistic

The theorem states that

$\exists N$ and $\exists w_j, \theta_j, \alpha_j$ for $j=1, \dots, N$ such that

$$|G(x) - f(x)| < \epsilon \quad \forall x \in I_n$$

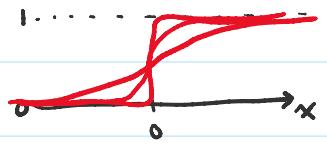


Proof:

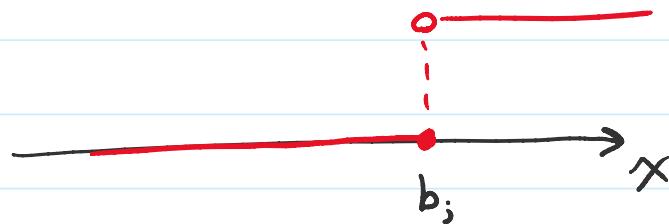
Since we let $w_{j,n} = \frac{\theta_j}{\alpha_j}$ for $j=1, \dots, N$

Suppose we let $w_j \rightarrow \infty$ for $j = 1, \dots, N$

Then $\sigma(w_j x) \xrightarrow{w_j \rightarrow \infty} \begin{cases} 0 & \text{for } x \leq 0 \\ 1 & \text{for } x > 0 \end{cases}$



Or $\sigma(w_j(x - b_j)) \xrightarrow{w_j \rightarrow \infty} \begin{cases} 0 & \text{for } x \leq b_j \\ 1 & \text{for } x > b_j \end{cases}$



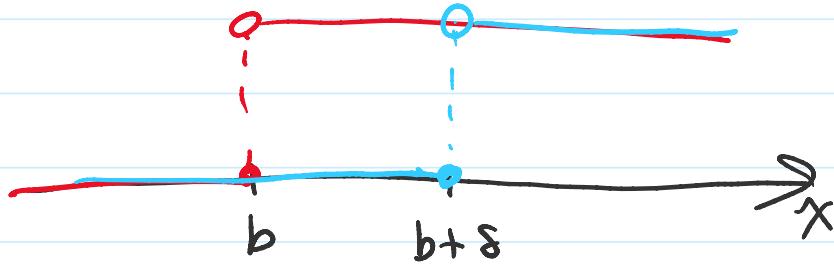
This is the same as the Heaviside step function,

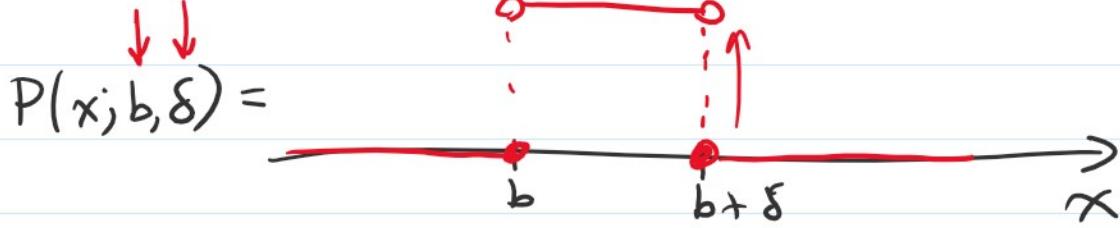
$$H(x) = \lim_{w_j \rightarrow \infty} \sigma(w_j x)$$

Define: $H(x; b) = \lim_{w \rightarrow \infty} \sigma(w(x - b))$

We can use two such functions to create a piece,

$$P(x; b, \delta) = H(x; b) - H(x; b + \delta)$$

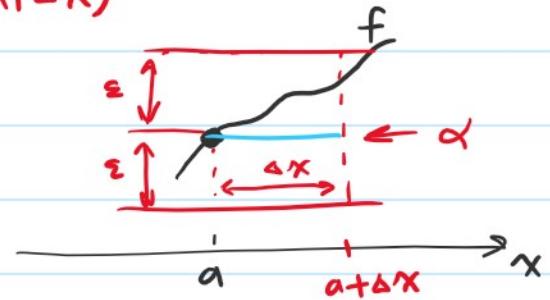




Since $f(x)$ is continuous, $\lim_{x \rightarrow a} f(x) = f(a)$. $\forall a \in I_n$

$\therefore \exists$ an interval, $(a, a+\Delta x)$ such that

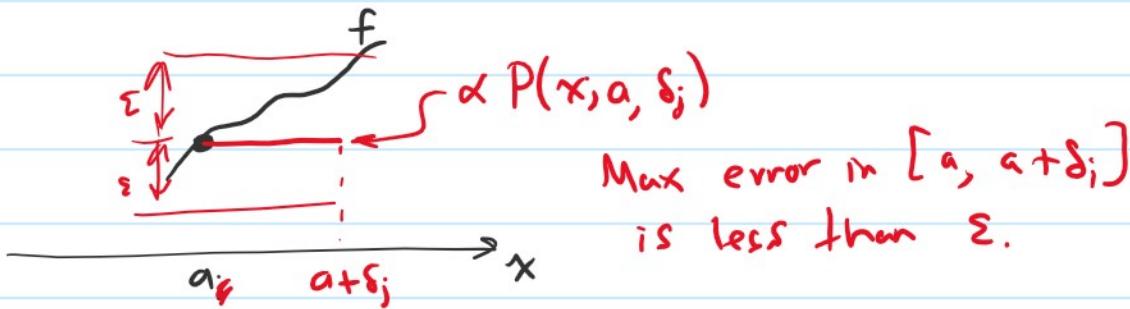
$$|f(x) - f(a)| < \varepsilon \quad \forall x \in (a, a+\Delta x)$$



Choose $b_i = a$, $\delta_i = \Delta x$, and $\alpha_i = f(a)$

$$\therefore |f(x) - f(a)| < \varepsilon \text{ for } a \leq x \leq a + \delta_i;$$

$$|f(x) - \alpha_i P(x; a, \delta_i)| < \varepsilon \quad ..$$



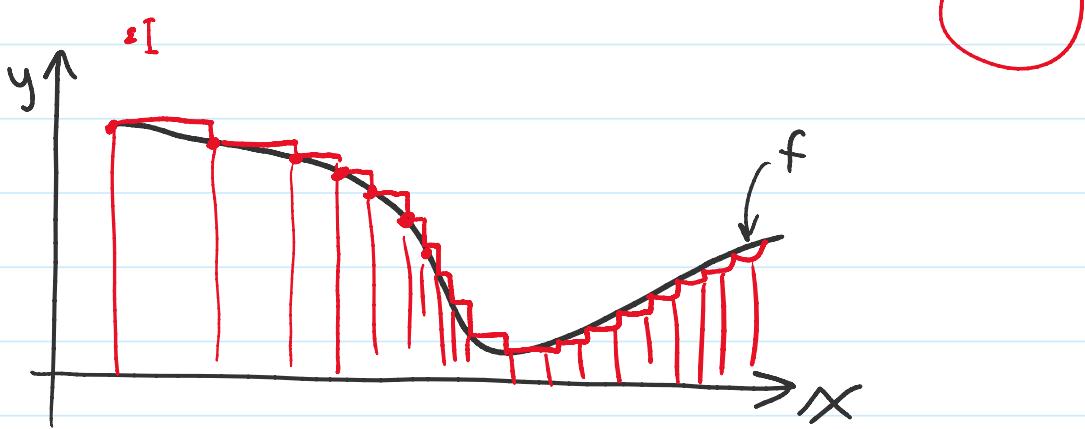
Repeat this process for $x = a = b_i + \delta_i$

Construct

$$G(x) = \sum_{j=1}^N \alpha_j P(x; b_j, \delta_j)$$

$\uparrow \varepsilon$





So, why would we ever need a neural network with more than one hidden layer?

Answer: The theorem guarantees existence,
but makes no claims about N , the
of hidden nodes.

N might grow exponentially as
 ϵ gets smaller.

Loss Functions

Goal: To become familiar with some of the most common ways to measure error.

We have to choose a way to quantify how close our output is to the target. For this, we use a "cost function", also known as an "objective function". There are many choices, but here are two commonly-used ones.

For input \vec{x} , our target is $\vec{t}(\vec{x})$, and the output of our network is $\vec{y}(\vec{x})$.

Mean Squared Error (MSE)

$$E(\vec{y}, \vec{t}) = \frac{1}{N} \sum_{i=1}^N \|\vec{y}_i - \vec{t}_i\|_2^2$$

N is the # of samples in the dataset.



The use of MSE as a cost function is often associated with linear activation functions, or ReLU. This is because these activation functions afford a larger output range.

Cross Entropy (Bernoulli Cross Entropy)

Consider a function (or network) with a single output that is either 0 or 1. The task of mapping inputs to the correct output (0 or 1) is a classification problem.

$$\vec{x} \rightarrow f(\vec{x}, \theta) \rightarrow y \in [0, 1]$$

where the true class is expressed in the target, $t \in \{0, 1\}$.

If we suppose that y is the probability that $x \rightarrow 1$,

$$\dots - P(y_i = 1 | x_i) = \text{[redacted]}$$

If we suppose that y is the probability that $x \rightarrow 1$,

$$y = P(x \rightarrow 1 | \theta) = f(x, \theta)$$

then we can treat it as a Bernoulli distribution.

i.e. $\begin{cases} P(x \rightarrow 1 | \theta) = y & \text{i.e. } t = 1 \\ P(x \rightarrow 0 | \theta) = 1 - y & \text{i.e. } t = 0 \end{cases}$

The likelihood of our data sample given our model is

$$P(x \rightarrow t | \theta) = y^t (1-y)^{1-t} \quad \log y^t (1-y)^{1-t} \\ = \log y^t + \log (1-y)^{1-t}$$

The task of "learning" would be finding a model (θ) that maximizes this likelihood.

Or, we could equivalently minimize the negative log-likelihood

$$E(y, t) = - (t \log y + (1-t) \log (1-y))$$

This log-likelihood formula is the basis of the ~~cross-entropy~~ cost function.
Given the dataset,

$$\{(x_1, t_1), (x_2, t_2), \dots, (x_N, t_N)\}$$

the expected cross entropy over the entire dataset is

$$-\mathbb{E}[t \log y + (1-t) \log (1-y)]_{\text{data}} \\ = -\frac{1}{N} \sum_{i=1}^N t_i \log y_i + (1-t_i) \log (1-y_i)$$

Cross entropy assumes that the output values are in the range [0, 1].

Hence, it works nicely with the logistic activation function.

Categorical Cross-Entropy (Multinoulli Cross-Entropy)

Consider a classification problem that has K classes ($K > 2$). Given an input,

the task of our model is to output the class of the input.

e.g. given an image of a digit, determine the digit class

Suppose the probability of getting class k is p_k .

$$\text{i.e. } p(k) = p_k \quad k=1, \dots, K, \quad \sum_k p_k = 1$$

Suppose we observe a sample of class k^*

Thus, the target class vector is

$$\vec{t} = [0 \ 0 \ \dots \underset{k^*}{1} \ \dots \ 0]^T$$

The likelihood of this sample is

$$P(\vec{t} | \vec{p}) = p_{k^*}$$

which can also be written

$$\begin{aligned} P(\vec{t} | \vec{p}) &= \prod_{k=1}^K p_k^{t_k} \\ &= p_1^{t_1} \cdot p_2^{t_2} \cdot \dots \cdot p_{k^*}^{t_{k^*}} \cdot \dots \cdot p_K^{t_K} \end{aligned}$$

Taking the logarithm of that gives

$$\log P(\vec{t} | \vec{p}) = \sum_{k=1}^K t_k \log p_k$$

Suppose our model is given the input \vec{x}

The network's output is

$$\vec{y} = f(\vec{x}, \theta) = P(\vec{t} | \theta, \vec{x})$$

We interpret \vec{y}_k as the estimated probability of \vec{x} being from class k

The true class is given by the target class vector \vec{t} .

Thus, the negative log-likelihood of \vec{x} is

$$E(y, t) = - \sum_{k=1}^K t_k \log y_k$$

The expected categorical cross-entropy for a dataset of N samples is

$$\mathbb{E}_{\text{data}} \left[- \sum_{k=1}^K t_k \log y_k \right] = \frac{-1}{N} \sum_{i=1}^N \sum_{k=1}^K t_{i,k} \log y_{i,k}$$

Since $\sum_k y_k = 1$, this cost function works well with

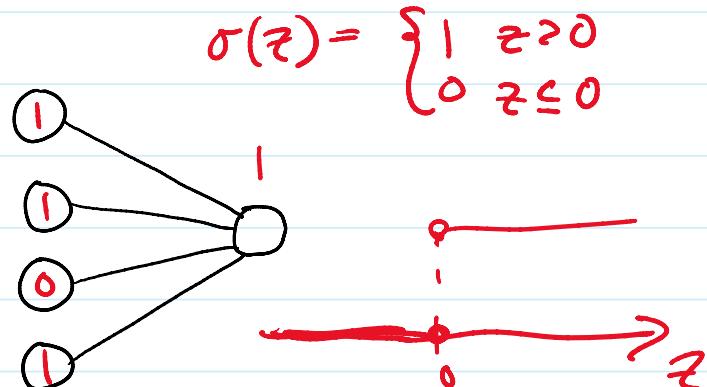
(Look at the new section at the end of Simple Neuron Models.)

Perceptrons

Goal: To see a simple neural learning algorithm, and understand its limitations.

Let's look at a simple neural network that aims to recognize certain input patterns.

For example, let the output node be a simple threshold neuron, and suppose we want the output node to be 1 when the input is $[1, 1, 0, 1]$, and zero otherwise.



Notice that if we set the weights to $[1, 1, 0, 1]$ (matching the input), then we maximize the input to the output node.

$$i) [1 \underset{\text{weights}}{1} \underset{\text{inputs}}{0} 1] \cdot [1 1 0 1] = 3 \rightarrow \sigma(3) = 1$$

But what about other inputs?

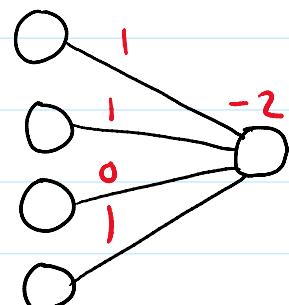
$$ii) [1 1 0 1] \cdot [0 1 1 0] = 1 \rightarrow \sigma(1) = 1$$

We need the un-matching input to give us a negative value so that the output node returns zero.

Solution: a negative bias

$$i) [1 1 0 1] \cdot [1 1 0 1] - 2 \rightarrow \sigma(-1) = 0$$

$$ii) [1 1 0 1] \cdot [0 1 1 0] - 2 \rightarrow \sigma(-1) = 0$$



Can we find the weights and bias automatically so that our perceptron produces the correct output for a variety of inputs?

To see an approach, let's look at a 2-D case.

~~Suppose the 4 different inputs are:~~

~~$[0, 0], [0, 1], [1, 0], \text{ and } [1, 1]$~~

~~And their corresponding outputs are~~

~~0 1 1 1 (an "OR" gate)~~

$$w = [0.6 \quad -0.2] \quad b = -0.1$$

~~And their corresponding outputs are~~

0, 1, 1, 1 (an "OR" gate)

Also, we will use the ~~E~~ error:

$$w = \begin{bmatrix} 0.6 \\ -0.2 \end{bmatrix} \quad b = -0.1$$

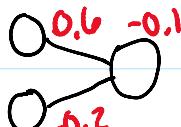
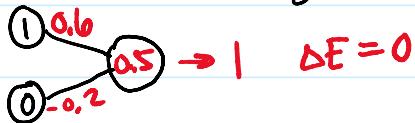
$$\Delta E = y - t$$

So, if $y=0, t=1 \rightarrow \Delta E = -1$

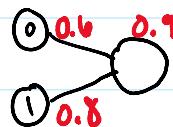
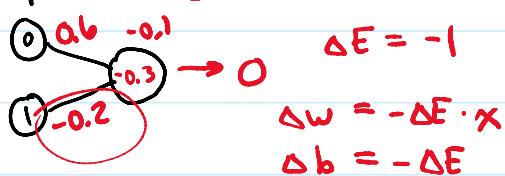
if $y=1, t=0 \rightarrow \Delta E = 1$

Start with random weights

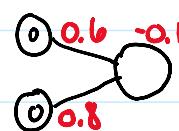
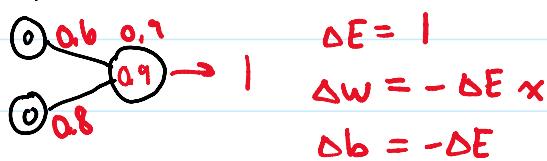
Input $[1, 0]$ target 1



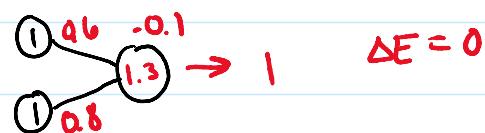
Input $[0, 1]$ target 1 $w \leftarrow w + \Delta w$



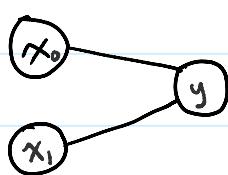
Input $[0, 0]$ target 0



Input $[1, 1]$ target 1



Graphical Interpretation

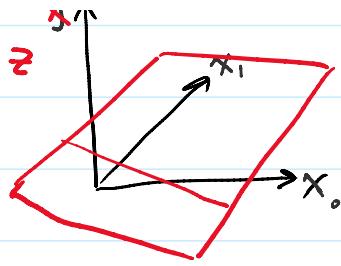


$$y = 0.6x_0 + 0.8x_1 - 0.1$$

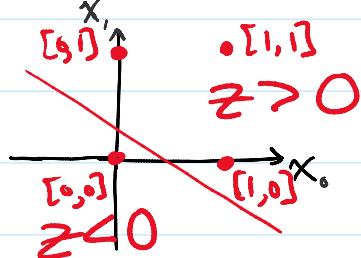
this is a linear equation... the equation of a plane in 3-D.



Looking down on the x_0 - x_1 plane,



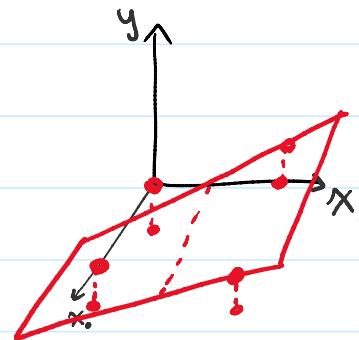
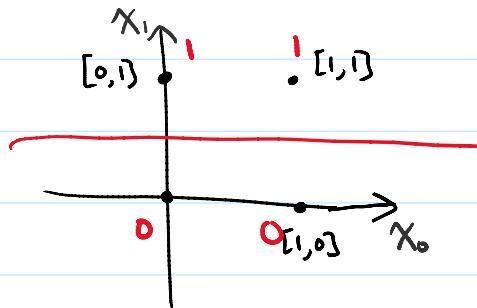
Looking down on the x_0 - x_1 plane,



Finding the weights and bias is the same as finding a linear classifier... a linear function that returns a positive value for the inputs that should yield a 1, and a negative value for the inputs that should yield a 0.

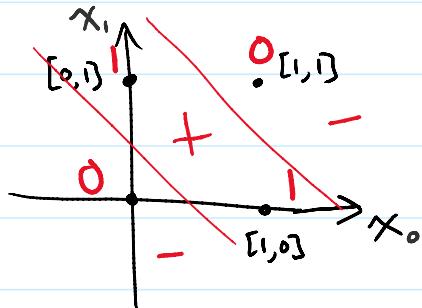
Another example:

$$[0, 0] \rightarrow 0 \quad [0, 1] \rightarrow 1 \quad [1, 0] \rightarrow 0 \quad [1, 1] \rightarrow 1 \quad \text{"echo } x_1\text{-value"}$$



A final example:

$$[0, 0] \rightarrow 0 \quad [0, 1] \rightarrow 1 \quad [1, 0] \rightarrow 1 \quad [1, 1] \rightarrow 0 \quad \text{"XOR"}$$



Perceptrons are simple, two-layer neural networks, so only work for linearly separable datasets.

If you want to handle non-linearly-separable data, your neural network is going to need more layers.

But they give us our first glimpse at a learning algorithm.

Gradient Descent Learning

Goal: To see how we can use a simple optimization method to tune our network weights.

The operation of our network can be written

$$\vec{y} = f(\vec{x}; \theta)$$

So, if our cost function is $E(\vec{y}, \vec{t})$, where \vec{t} is the target, then neural learning becomes the optimization problem

$$E = \min_{\theta} E \left[E(f(\vec{x}; \theta), \vec{t}(\vec{x})) \right]_{\text{data}}$$

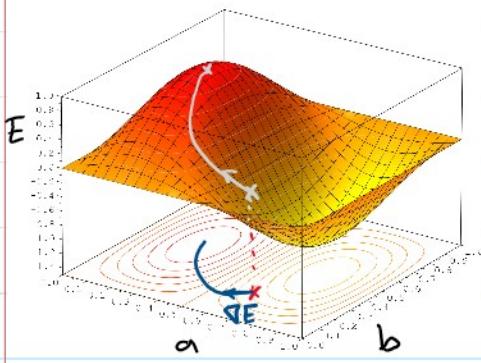
We can apply gradient descent to E , using the gradient

$$\nabla_{\theta} E = \left[\frac{\partial E}{\partial \theta_0}, \frac{\partial E}{\partial \theta_1}, \dots, \frac{\partial E}{\partial \theta_p} \right]^T$$

Gradient-Based Optimization

If you want to find a local maximum of a function, you can simply start somewhere, and keep walking uphill.

For example, suppose you have a function with two inputs, $E(a, b)$. You wish to find a and b to maximize E .



[https://commons.wikimedia.org/wiki/File:2D_Wavefunction_\(2.1\)_Surface_Plot.png](https://commons.wikimedia.org/wiki/File:2D_Wavefunction_(2.1)_Surface_Plot.png)

We are trying to find the parameters (\bar{a}, \bar{b}) that yield the maximum value of E .

$$\text{i.e. } (\bar{a}, \bar{b}) = \arg \max_{(a, b)} E(a, b)$$

No matter where you are, "uphill" is in the direction of the gradient vector,

$$\nabla E(a, b) = \left[\frac{\partial E}{\partial a}, \frac{\partial E}{\partial b} \right]^T$$

Gradient ascent is an optimization method where you step in the direction of your gradient vector.

If your current position is (a_n, b_n) , then

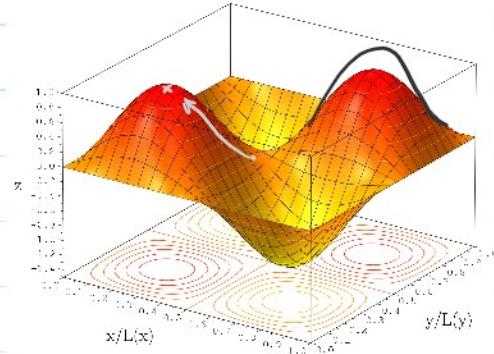
$$(a_{n+1}, b_{n+1}) = (a_n, b_n) + k \nabla E(a_n, b_n)$$

where k is your step multiplier.

where k is your step multiplier.

Gradient **DESCENT** aims to **minimize** your objective function. So, you walk downhill, stepping in the direction **opposite** the gradient vector.

Note that there is no guarantee that you will actually find the global optimum. In general, you will find a local optimum that may or may not be the global optimum.



Approximating the Gradient Numerically

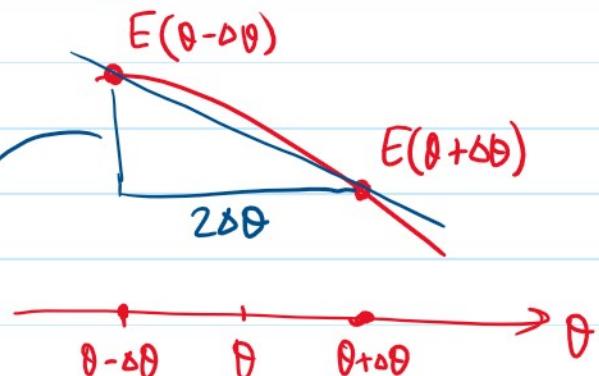
We can estimate the partial derivatives in the gradient using finite-differencing.

Finite-Difference Approximation

For a function $E(\theta)$, we can approximate

$\frac{\partial E}{\partial \theta}$ using

$$\frac{dE}{d\theta} \approx \frac{E(\theta + \Delta\theta) - E(\theta - \Delta\theta)}{2\Delta\theta}$$



As an example, consider this network:



It's a neural network, with connection weights and biases shown.

We formulate the problem as

$$\min_{\theta} E(f(\vec{x}, \theta), \vec{f}(\vec{x}))$$

Or, more compactly, $\min_{\theta} \bar{E}(\theta)$ where $\bar{E}(\theta) = E(f(\vec{x}, \theta), \vec{f}(\vec{x})) = (y - t)^2$

Consider θ_1 on its own. With $\theta_1 = -0.01$, our network output is

$$y = 0.5089$$

This gives

$$\bar{E}(0.01) = 0.24113$$

What if we perturb θ_1 , so that $\theta_1 = -0.01 + 0.5 = 0.49$.

Then our output is

$$y = 0.5693$$

This yields

$$\bar{E}(0.49) = 0.240761$$

If, instead, we perturb θ_1 so that $\theta_1 = -0.01 - 0.5 = -0.51$,

then our output is

$$y = 0.5086$$

which gives

$$\bar{E}(-0.51) = 0.241509$$

In summary,

Parameters	MSE
$\theta_1 + \Delta\theta$	0.240761
θ_1	0.24113
$\theta_1 - \Delta\theta$	0.241509

We can estimate $\frac{\partial \bar{E}}{\partial \theta_1}$ using

$$\frac{\partial \bar{E}}{\partial \theta_1} \approx \frac{\bar{E}(0.49) - \bar{E}(-0.51)}{2(0.5)} = -0.0007475$$

Obviously, increasing θ_1 seems to be the right thing to do.

$$\theta_1 \leftarrow -0.01 - K(-0.0007475)$$

$$= -0.01 + K(0.0007475)$$

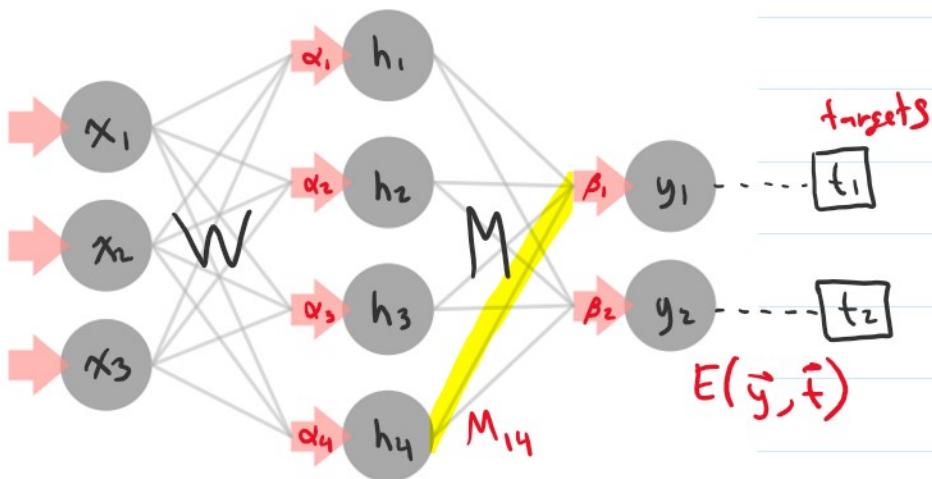
(demo: XOR example)

Error Backpropagation

Goal: To find an efficient method to compute the gradients for gradient-descent optimization.

We can apply gradient descent on a multi-layer network, ~~again~~ using chain rule to calculate the gradients of the error with respect to deeper connection weights and biases.

Consider the network



α_i is the input current to hidden node i .

β_j is the input current to output node j .

For our cost (loss) function, we will use $E(\bar{y}, \bar{t})$

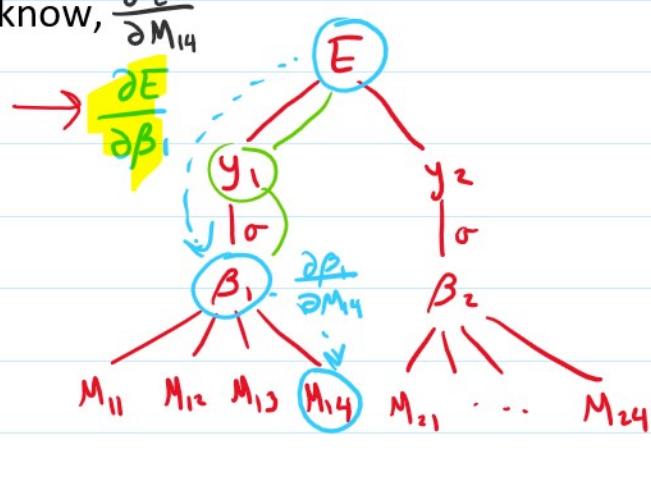
For learning, suppose we want to know, $\frac{\partial E}{\partial M_{14}}$

e.g. M_{14}

$$\frac{\partial E}{\partial M_{14}} = \frac{\partial E}{\partial \beta_1} \frac{\partial \beta_1}{\partial M_{14}}$$

Recall, $E(\bar{y}, \bar{t})$

$$= E\left(\sigma\left(\bar{M}\bar{h} + \bar{b}\right), \bar{t}\right)$$



$$= E \left(\underbrace{\sigma(Mh + b)}_{\vec{\beta}}, \tau \right)$$

$$\therefore \frac{\partial E}{\partial \beta_1} = \frac{\partial E}{\partial y_1} \frac{dy_1}{d\beta_1} \quad \leftarrow \text{we'll return to this later}$$

Thus,

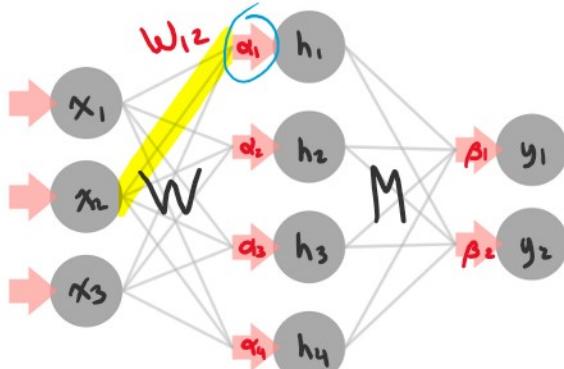
$$\frac{\partial E}{\partial M_{14}} = \frac{\partial E}{\partial y_1} \frac{dy_1}{d\beta_1} \frac{\partial \beta_1}{\partial M_{14}}$$

$$\text{Recall } \beta_1 = \sum_{i=1}^4 M_{1i} h_i + b_1 \\ = M_{11} h_1 + M_{12} h_2 + M_{13} h_3 + M_{14} h_4 + b_1$$

$$\therefore \frac{\partial E}{\partial M_{14}} = \boxed{\frac{\partial E}{\partial y_1} \frac{dy_1}{d\beta_1}} h_4$$

$$\frac{\partial \beta_1}{\partial M_4} = h_4$$

OK, that works for the connection weights between the top two layers. What about the connection weights between layers deeper in the network? $e.g.$ W_{12} of F



$$\frac{\partial E}{\partial w_{12}} = \frac{\partial E}{\partial \alpha_1} \frac{\partial \alpha_1}{\partial w_{12}}$$

$$\frac{\partial E}{\partial x_1} = \frac{\partial E}{\partial h_1} \frac{dh_1}{dx_1}$$

$$= \left(\frac{\partial E}{\partial \beta_1} \frac{\partial \beta_1}{\partial h_i} + \frac{\partial E}{\partial \beta_2} \frac{\partial \beta_2}{\partial h_i} \right) \frac{dh_i}{1}$$

$$\alpha_i = \sum_{j=1}^n w_{ij} x_j + a_j$$

$$\therefore \frac{\partial \alpha_1}{\partial w_{12}} = x_2$$

$$= \left(\frac{\partial E}{\partial \beta_1} \frac{\partial \beta_1}{\partial h_i} + \frac{\partial E}{\partial \beta_2} \frac{\partial \beta_2}{\partial h_i} \right) \frac{dh_i}{d\alpha_i}$$

$$= \left(\frac{\partial E}{\partial \beta_1} M_{1i} + \frac{\partial E}{\partial \beta_2} M_{2i} \right) \frac{dh_i}{d\alpha_i}$$

We computed these already when we were learning M.

$$= (M_{1i}, M_{2i}) \cdot \left(\frac{\partial E}{\partial \beta_1}, \frac{\partial E}{\partial \beta_2} \right) \frac{dh_i}{d\alpha_i}$$

More generally, $\vec{x} \in \mathbb{R}^X$, $\vec{h} \in \mathbb{R}^H$, $\vec{y}, \vec{t} \in \mathbb{R}^Y$

$$\frac{\partial E}{\partial \alpha_i} = \underbrace{[M_{1i} \dots M_{Yi}]}_{\text{This is the } i^{\text{th}} \text{ column of } M} \cdot \begin{bmatrix} \frac{\partial E}{\partial \beta_1} & \dots & \frac{\partial E}{\partial \beta_Y} \end{bmatrix} \frac{dh_i}{d\alpha_i}$$

This is the i^{th} column of M

$$= [M_{1i} \dots M_{Yi}] \begin{bmatrix} \frac{\partial E}{\partial \beta_1} \\ \vdots \\ \frac{\partial E}{\partial \beta_Y} \end{bmatrix} \frac{dh_i}{d\alpha_i}$$

For more elements

$$\begin{bmatrix} \frac{\partial E}{\partial \alpha_1} \\ \frac{\partial E}{\partial \alpha_2} \\ \vdots \\ \frac{\partial E}{\partial \alpha_n} \end{bmatrix} =$$

For all elements,

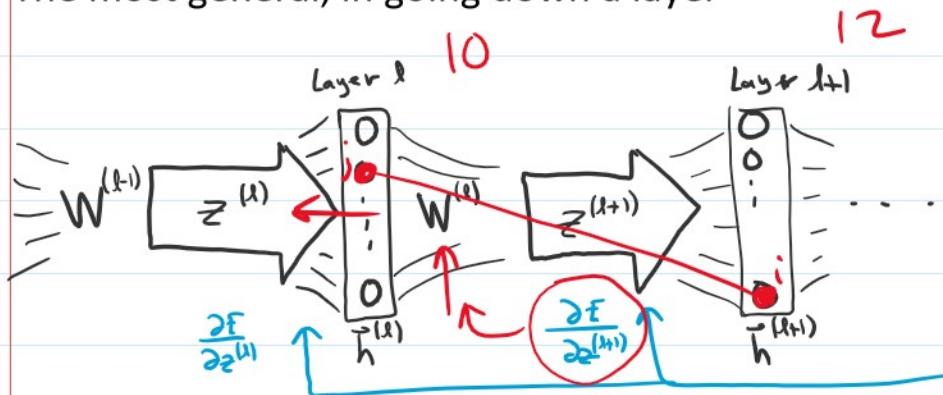
$$\begin{bmatrix} \frac{\partial E}{\partial \alpha_1} \\ \vdots \\ \frac{\partial E}{\partial \alpha_n} \end{bmatrix} = \begin{bmatrix} M_{11} & \dots & M_{1H} \\ \vdots & & \vdots \\ M_{Y1} & \dots & M_{YH} \end{bmatrix} \begin{bmatrix} \frac{\partial E}{\partial \beta_1} \\ \vdots \\ \frac{\partial E}{\partial \beta_Y} \end{bmatrix} \circ \begin{bmatrix} \frac{dh_1}{d\alpha_1} \\ \vdots \\ \frac{dh_H}{d\alpha_n} \end{bmatrix}$$

Hadamard Product
 $[a b] \odot [c d] = [ac bd]$

$$\nabla_\alpha E = \frac{\partial E}{\partial \alpha} = \frac{d\vec{h}}{d\vec{\alpha}} \odot M^T \frac{\partial E}{\partial \beta} = \frac{d\vec{h}}{d\vec{\alpha}} \odot M^T \nabla_\beta E$$

The most general, in going down a layer

The most general, in going down a layer



Suppose we have $\frac{\partial E}{\partial \bar{z}^{(l+1)}} = \nabla_{\bar{z}^{(l+1)}} E$

Let $\bar{h}^{(l+1)} = \sigma(\bar{z}^{(l+1)}) = \sigma(W^{(l)} \bar{h}^{(l)} + b^{(l+1)})$

$$\frac{\partial E}{\partial \bar{z}^{(l)}} = \frac{d\bar{h}^{(l)}}{d\bar{z}^{(l)}} \Theta \left[W^{(l)} \right]^T \frac{\partial E}{\partial \bar{z}^{(l+1)}}$$

Then, to compute $\frac{\partial E}{\partial w_{ij}^{(l)}}$...

$$\frac{\partial E}{\partial w_{ij}^{(l)}} = \frac{\partial E}{\partial \bar{z}_i^{(l+1)}} \frac{\partial \bar{z}_i^{(l+1)}}{\partial w_{ij}^{(l)}}$$

$$= \frac{\partial E}{\partial \bar{z}_i^{(l+1)}} h_j^{(l)}$$

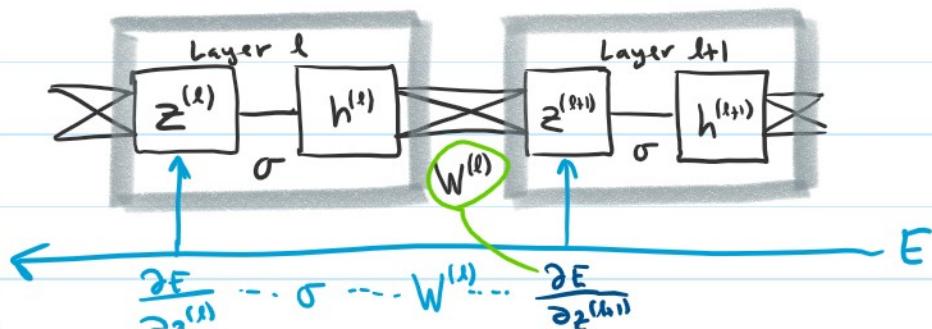
$$\frac{\partial E}{\partial w^{(l)}} = \begin{bmatrix} 1 \\ \frac{\partial E}{\partial \bar{z}^{(l+1)}} \\ 1 \end{bmatrix} \left[-\bar{h}^{(l)} \right] =$$

$$\bar{z}^{(l+1)} = W^{(l)} \bar{h}^{(l)} + \dots$$

Same size
as
 $w^{(l)}$

outer product

Another graphical summary



$$\frac{\partial E}{\partial z^{(1)}} \dots \sigma \dots W^{(1)} \dots \frac{\partial E}{\partial z^{(L+1)}}$$

$$\frac{\partial E}{\partial z^{(1)}} = \tilde{\sigma}'(\tilde{z}^{(1)}) \odot (\tilde{W}^{(1)})^T \frac{\partial E}{\partial z^{(L+1)}}$$

$$\frac{\partial E}{\partial W^{(1)}} = \frac{\partial E}{\partial z^{(L+1)}} [h^{(1)}]^T \quad (\text{outer prod}).$$

END

Training and Testing

Goal: Develop a process to use our labelled data to generate models that can predict future, unseen samples.

We have seen how to adjust a neural network to get it to learn our training data. Of course, the purpose of training a model is so that we can use it on other samples that aren't in our training set. For this reason, we usually break our data into two pieces:

1. **Training set:** Use most of your labeled data to train your model.
2. **Test set:** Once your model is trained, use the remaining labeled samples to evaluate your model.

Why? Suppose after some trial-and-error, adjusting the hyperparameters, such as:

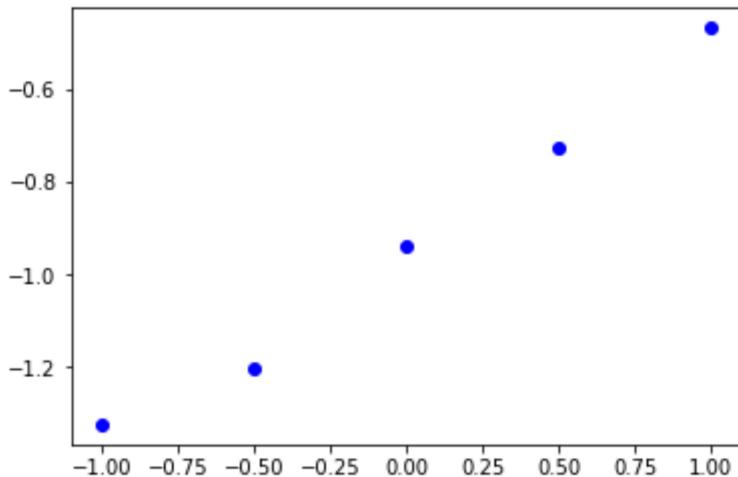
- number of neurons in each layer
- learning rate
- number of epochs
- initial weights

you finally get a low error on the training data. Does that accomplish what you want? Consider an example.

As an experiment, consider noisy samples coming from the ideal mapping,

$$y = 0.4x - 0.9$$

We can only get noisy samples from that mapping.



Our training dataset has 5 samples. And since this is a regression problem, we will use a **identity** activation function on the output, and **MSE** as a loss function.

Training usually entails going through the training data repeatedly, updating the network weights as we go. Each pass through the data is called an **epoch**

Let's create a neural network to learn this mapping.

```
net = Network([1, 1000, 1])
    ↑ # inputs      ↑ # hidden      ↑ # outputs.
```

Before training...

Training MSE = 0.955758517256

Call the learning function for many epochs.

```
net.Learn([training_input, training_output], epochs=40000, lrate=0.01)
```

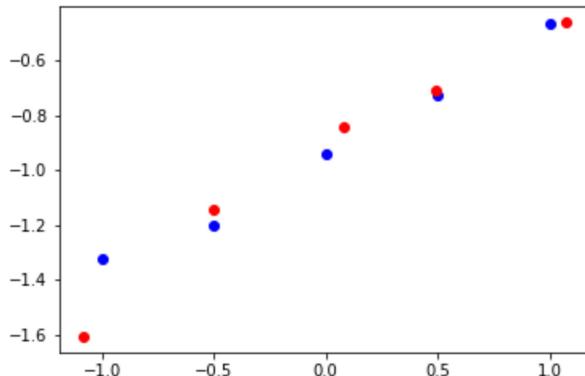
After going through the dataset MANY times, our average loss is

Training MSE = 0.000692658082161

SUCCESS!

Let's bask in the glory of our brilliance, and demonstrate how

great we are on another sampling of data.



Red samples
are 5 new
observations

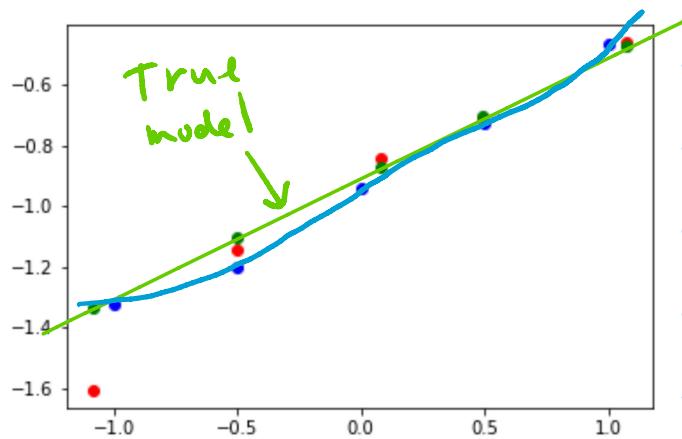
Our average loss on this new set of samples is

Test MSE = 0.0156456136501

Uh-oh!

This is not as good as our training error!

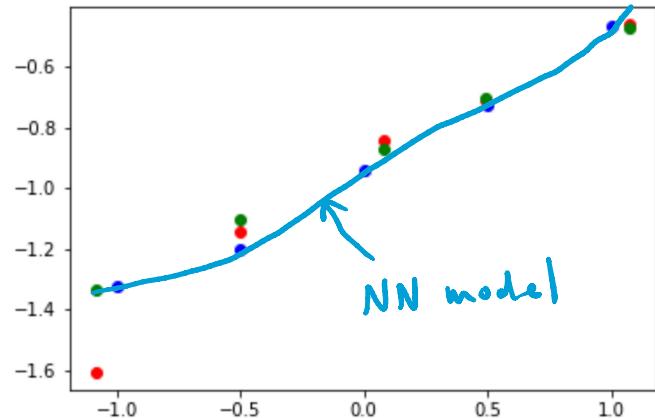
In fact, what happens when we give it a perfect dataset, without noise?



Perfect MSE = 0.00161413205859

Even perfect data has a higher error than the test data. ???

The false sense of success we get from the results on our training dataset is known as **overfitting**.



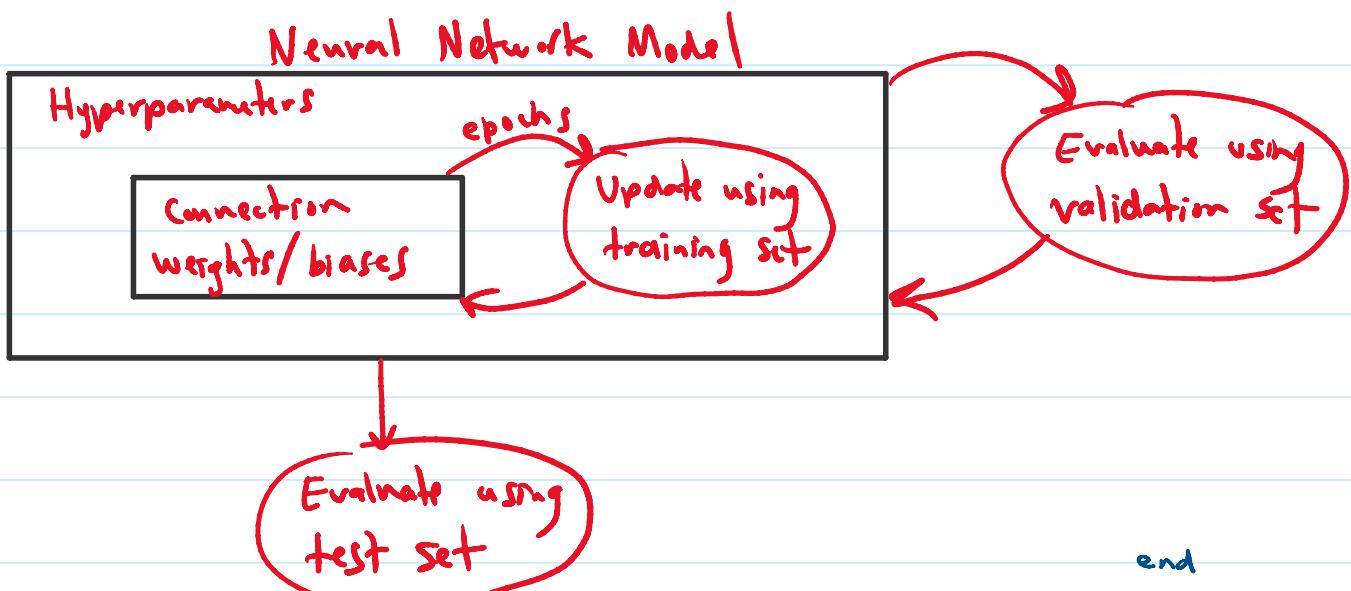
The model starts to fit the noise specific to the training set, rather than just to the underlying model.

Recall that our sole purpose was to create a model to predict the output for samples it hasn't seen. How can we avoid overfitting?

Validation

If we want to estimate how well our model will generalize to samples it hasn't trained on, we can withhold part of the training set and try our model on that "validation set". Once our model does reasonably well on the validation set, then we have more confidence that it will perform reasonably well on the test set.

It's common to use a random subset of the training set as a validation set.



(" ")
test set

end

Overfitting

Goal: See some tricks for how to mitigate ~~against~~ overtraining.

We saw that if a model has enough degrees of freedom, it can become hyper-adapted to the training set, and start to fit the noise in the dataset.

Training error is very small

This is a problem because the model does not generalize well to new samples.

Test error is much bigger than training error

There are some strategies to try to stop our network from trying to fit the noise.

Regularization

Weight Decay

We can limit overfitting by creating a preference for solutions with smaller weights, achieved by adding a term to the loss function that penalizes for the magnitude of the weights.

$$\tilde{E}(y, t; \theta) = E(y, t; \theta) + \lambda \|\theta\|_F^2$$

$\|\theta\|_F = \sqrt{\sum_j \theta_j^2}$ "Frabinius Norm"

$$\|\theta\|_F^2 = \theta_1^2 + \theta_2^2 + \dots + \theta_n^2$$

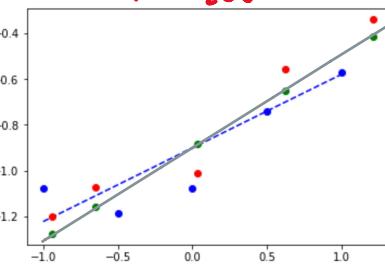
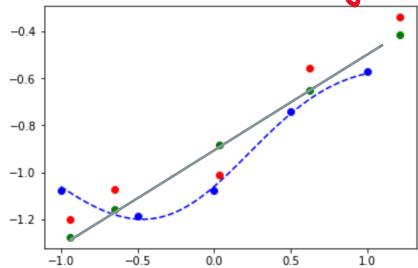
How does this change our gradients, and thus our update rule?

$$\frac{\partial \tilde{E}}{\partial \theta_i} = \frac{\partial E}{\partial \theta_i} + 2\lambda \theta_i$$

$$\theta_{\text{new}} \leftarrow \theta_{\text{old}} - k \frac{\partial \tilde{E}}{\partial \theta_{\text{old}}} \quad \text{decay}$$

$\lambda = 0$ (no reg'n)

$\lambda = \frac{1}{200}$



$$\theta_{\text{old}} - k \frac{\partial E}{\partial \theta_{\text{old}}} - \boxed{2k\lambda \theta_{\text{old}}}$$

Training MSE = 0.00596061335994
 Test MSE = 0.00574788521961
 Perfect MSE = 0.00178339752734

λ controls the weight of the regularization term.

One can also use different norms. For example, it is common to use the L1 norm,

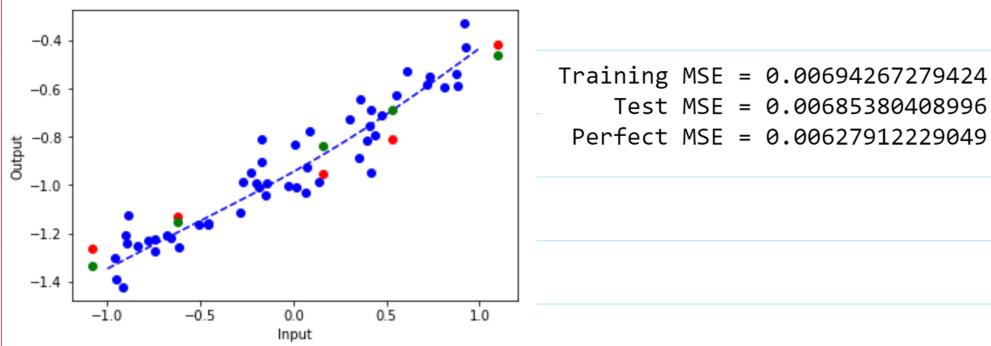
$$L_1(\theta) = \sum_i |\theta_i|$$

$$\frac{\partial L_1}{\partial \theta} = \text{sign}(\theta)$$

The L1 norm tends to favour sparsity (most weights are close to zero, with only a small number of non-zero weights).

Data Augmentation

Another approach is to include a wider variety of samples in your training set, so that the model is less likely to focus its efforts on the noise of a few.



Where does this extra data come from?

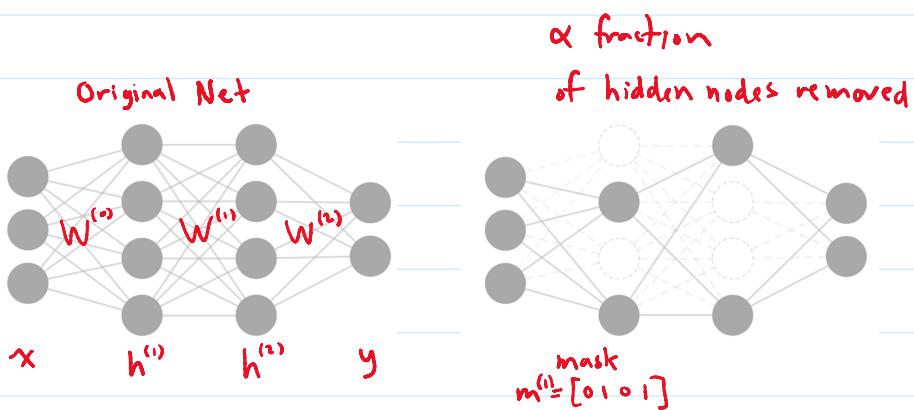
For image-recognition datasets, one can generate more samples by shifting or rotating the images. Those transformations presumably do not change the labelling.

Dropout

The last method we will talk about is the most bizarre.

While training using the dropout method, you systematically ignore a large fraction (typically half) of the hidden nodes for each sample. That is, given a dropout probability, α , each hidden node will be dropped with probability α .

A dropped node is temporarily taken off-line and set to zero.



Do both a feedforward and backprop pass with this diminished network.

$$z^{(1)} = W^{(1)}x + b^{(1)} \Rightarrow \hat{h}^{(1)} = \sigma(z^{(1)}) \in \mathbb{R}^{N^{(1)}}$$

Then apply mask $m^{(1)} = \{0, 1\}^{N^{(1)}}$

$$h^{(1)} = \hat{h}^{(1)} \odot m^{(1)}$$

Important caveat:

Removing a bunch of the hidden nodes reduces the input current to the next layer. The weights learned in this context will not work properly in the full network.

We have to **scale** the activities from a diminished layer in order to give **reasonable** inputs to the next layer.

We scale up the remaining activities of the diminished layer.

Suppose only R of the $N^{(1)}$ nodes remained.

Let $\bar{h}^{(1)}$ be the scaled activities, where $\bar{h}^{(1)} = \frac{N^{(1)}}{R} h^{(1)}$

This give us

$$\mathbb{E}\left[\sum_i \bar{h}_i^{(1)}\right] = \mathbb{E}\left[\sum_i \hat{h}_i^{(1)}\right]$$

↑
on average
 $\frac{N^{(1)}}{R} = \frac{1}{1-\alpha}$

We store these scaled activities in our hidden nodes, but also record the scale factor, $\frac{N^{(1)}}{R}$.

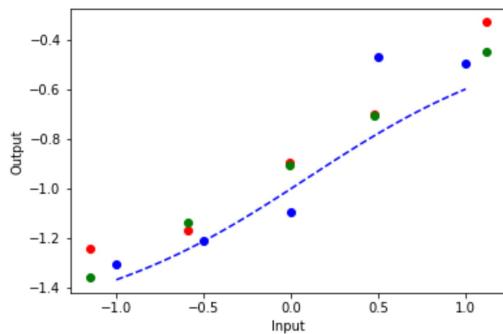
During backprop, suppose we have to make sure we undo this scaling as we pass the gradient back. We want to apply weight updates as though we were using the full network (no dropouts).

For example,

$$\frac{\partial E}{\partial W^{(1)}} = \underbrace{\frac{\partial E}{\partial z^{(2)}}}_{\downarrow} \underbrace{\frac{\partial z^{(2)}}{\partial W^{(1)}}}_{\downarrow} = \frac{\partial E}{\partial z^{(2)}} \bar{h}^{(1)} = \frac{\partial E}{\partial z^{(2)}} \left(\frac{R}{N^{(1)}} h^{(1)} \right)$$

We have to un-scale the hidden node activities during backprop

We have to un-scale the hidden node activities during backprop.



Training MSE = 0.0234801932225
Test MSE = 0.0212572143563
Perfect MSE = 0.00874342134987

Why does dropout work?

- It's akin to training a bunch of different networks and combining their answers. Each diminished network is like a contributor to this consensus strategy.
- Dropout disallows sensitivity to particular combinations of nodes. Instead, the network has to seek a solution that is robust to loss of nodes.

End of L12

Enhancing Optimization

Goal: To learn some methods that help learning go faster.

Suppose our training set is

$$\{(x_1, t_1), \dots, (x_p, t_p)\} \text{ where } \vec{x} \in \mathbb{R}^X, \vec{t} \in \mathbb{R}^Y$$

Notice that we can put all of our inputs into a single matrix

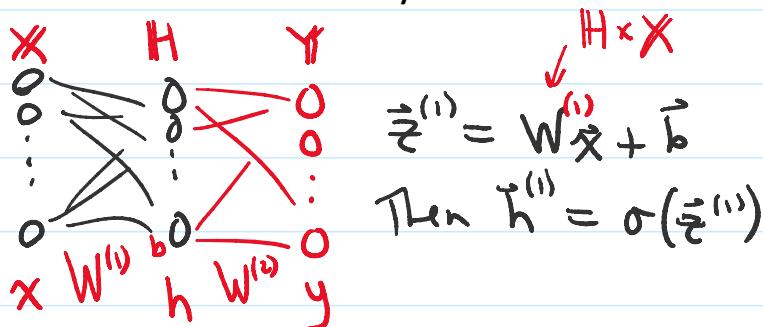
$$X = [x_1 | x_2 | \dots | x_p]$$

As well as our targets

$$T = [t_1 | \dots | t_p]$$

Does this help us? Yes!

Consider the 1st hidden layer...



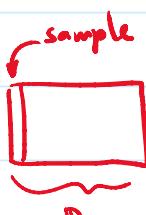
But we can process all P inputs at once!

$$\vec{z}^{(1)} = W^{(1)}X + b \underbrace{[1 \dots 1]}_P$$

$$\boxed{\vec{x}} = \boxed{b} | \boxed{b} | \dots | \boxed{b}$$

$$\text{Then, } \vec{H}^{(1)} = \sigma(\vec{z}^{(1)})$$

$$\text{At the top layer, we get } Y \in \mathbb{R}^{Y \times P} = \boxed{y}$$



$$E(Y, T) = \frac{1}{P} \sum_{p=1}^P E(y_p, t_p)$$

Now, working our way back down,

$$\partial E \quad \vdots \rightarrow \perp \quad \nwarrow \circ$$

Now, working our way back down,

$$\frac{\partial E}{\partial z^{(4)}} = (\gamma - T) \frac{1}{P} \quad Y \times P$$

And going down one layer

$$\frac{\partial E}{\partial z^{(3)}} = \frac{dH}{dz^{(3)}} \odot (W^{(2)})^T \frac{\partial E}{\partial z^{(4)}}$$

$H \times P$ $H \times P$ $H \times Y$ $Y \times P$

Then,

$$\frac{\partial E}{\partial W^{(2)}} = \frac{\partial E}{\partial z^{(3)}} (H^{(2)})^T$$

$Y \times H$ $Y \times P$ $P \times H$

$H=20 \quad Y=10$

$$= \sum_{i=1}^{10} \left[\begin{array}{c} \uparrow \\ \frac{\partial E}{\partial z_i^{(3)}} \\ \downarrow \end{array} \right] \left[\begin{array}{c} \leftarrow h_i^{(2)} \rightarrow \\ \dots \\ \leftarrow h_p^{(2)} \rightarrow \end{array} \right] + \dots + \left[\begin{array}{c} \uparrow \\ \frac{\partial E}{\partial z_{20}^{(3)}} \\ \downarrow \end{array} \right] \left[\begin{array}{c} \leftarrow h_i^{(2)} \rightarrow \\ \dots \\ \leftarrow h_p^{(2)} \rightarrow \end{array} \right]$$

$$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \begin{bmatrix} 4 & 5 \\ 8 & 10 \\ 12 & 15 \end{bmatrix}$$

Hence, we can use the same formulas to process a whole batch of samples.

One problem: processing the entire dataset for a single update to the weights can be slow. Instead, there is an intermediate approach.

Stochastic Gradient Descent

Computing the gradient of the cost function can be very expensive and time-consuming, especially if you have a huge training set.

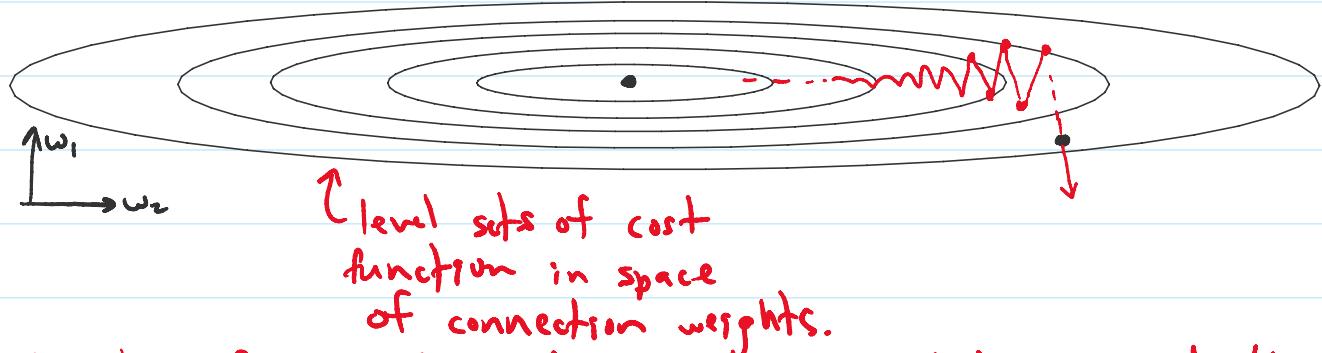
Rather than compute the full gradient, we can try to get a cheaper estimate by computing the gradient from a random sampling.

We use the estimate from this batch to update our weights, and then choose subsequent batches from the remaining samples.

This method is called Stochastic Gradient Descent.

Momentum

Consider gradient descent optimization in this situation...



The shape of the cost function causes oscillation, and this back-and-forth action can be inefficient.
Instead, we can smooth out our trajectory using
Recall from physics,

$$\frac{d\mathbf{D}}{dt} = \mathbf{V}$$

(velocity)

$$\frac{d\mathbf{V}}{dt} = \mathbf{A} - r\mathbf{V}$$

(acceleration)

So, solving numerically using Euler's method...

$$\mathbf{D}_{n+1} = \mathbf{D}_n + \Delta t \mathbf{V}_n \quad \textcircled{1}$$

$$\mathbf{V}_{n+1} = (1-r) \mathbf{V}_n + \Delta t \mathbf{A}_n$$

r is resistance from friction

when driving, this is what you control.
e.g. gas pedal, brakes, steering.

In our previous optimization method, \mathbf{D} represented our weights, and \mathbf{V} was like our error gradients. Then we updated our weights using $\textcircled{1}$

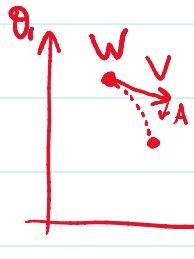
$$\text{Distance} \quad \mathbf{D}_{n+1} = \mathbf{D}_n + \Delta t \mathbf{V}_n$$

$$\text{Weights} \quad \mathbf{W}_{n+1} = \mathbf{W}_n - K \frac{\partial E}{\partial \mathbf{W}_n}$$

But we can instead treat our error gradients as \mathbf{A} , and integrate to get an accumulated weight update akin to \mathbf{V}

But we can instead treat our error gradients as \mathbf{A} , and integrate to get an accumulated weight update, akin to \mathbf{V} .
In fact, let's call it \mathbf{V} .

It's like our weights are dictated by our location in parameter space,



and we move around weight space, accelerated by the error gradients. We build speed if we get a lot of acceleration in the same direction.

We gain momentum .

For each weight W_{ij} , we also calculate V_{ij} .
Or, in matrix form, for each $W^{(l)}$, we have $V^{(l)}$

$$V^{(l)} \leftarrow (1-r)V^{(l)} + \frac{\partial E}{\partial W^{(l)}}$$

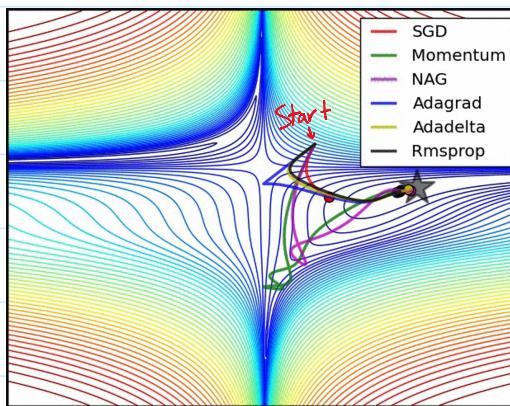
Or, as is commonly used,

$$V^{(l)} \leftarrow \beta V^{(l)} + (1-\beta) \frac{\partial E}{\partial W^{(l)}}$$

Then, update our weights using

$$W^{(l)} \leftarrow W^{(l)} - KV^{(l)}$$

Not only does this smooth out oscillations, but can also help to avoid getting stuck in local minima.



<http://ruder.io/optimizing-gradient-descent/>

End of 14th lecture

Deep Neural Networks

Goal: To see the advantages and disadvantages of deep neural networks: representational power vs. vanishing or exploding gradients.

How many layers should our neural network have?

Recall the Universal Approximation Theorem:

Theorem 2. *Let σ be any continuous sigmoidal function. Then finite sums of the form*

$$G(x) = \sum_{j=1}^N \alpha_j \sigma(\omega_j x + \theta_j)$$

are dense in $C(I_n)$. In other words, given any $f \in C(I_n)$ and $\varepsilon > 0$, there is a sum, $G(x)$, of the above form, for which

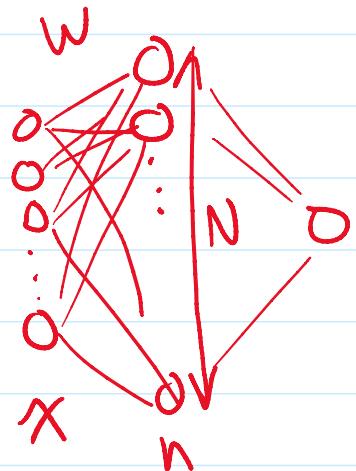
$$|G(x) - f(x)| < \varepsilon \quad \text{for all } x \in I_n.$$

Cybenko G, "Approximation by Superpositions of a Sigmoidal Function", *Math. Control Signals Systems*, 2:303-314, 1989.

Thus, we really only ever need one hidden layer. But is that the best approach, either in number of nodes, or learning efficiency? No, it can be shown that such a shallow network could require an exponentially large number of nodes (ie. A really big N) to work.

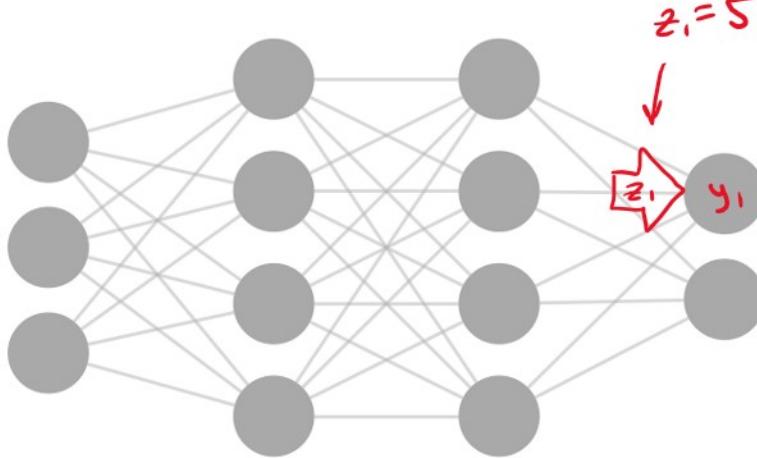
So, a deeper network is preferred in many cases.
(see Python example)

So, why don't we always use really deep networks?



Vanishing Gradients

Suppose the initial weights and biases were large enough that the input current to many of the nodes was not too close to zero. As an example, consider one of the output nodes.



$$y_1 = \sigma(z_1)$$

$$= \frac{1}{1+e^{-5}}$$

$$= 0.9933\dots$$

$$\frac{dy_1}{dz_1} = y_1(1-y_1)$$

$$= 0.0066$$

Compare that to if the input current was 0.1.

$$y_1 = \sigma(0.1) = 0.525$$

$$\frac{dy}{dz_1} = 0.249 \text{ Almost 40x larger than}$$



Hence, the updates to the weights will be smaller when the input currents are large in magnitude.

What about the next layer down?

$$\text{Suppose } \frac{\partial E}{\partial z^{(4)}} \approx 0.1$$

What if the inputs to the penultimate layer were around 4 in magnitude?

Then the corresponding slopes of their sigmoid functions will also be small.

$$\sigma(4) = 0.982$$

$$\sigma'(4) = 0.0177$$

Recall that

$$\frac{\partial E}{\partial z^{(3)}} = \frac{d h^{(2)}}{d z^{(3)}} \odot \left(W^{(3)} \right)^T \frac{\partial E}{\partial z^{(4)}}$$

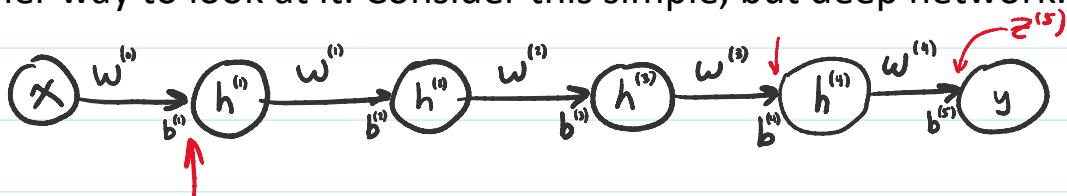
$$\begin{aligned}\frac{\partial E}{\partial z^{(3)}} &= \frac{dh^{(4)}}{dz^{(3)}} \odot (W^{(3)})^T \frac{\partial E}{\partial z^{(4)}} \\ &= (\text{approx } 0.0177) \odot (W^{(3)})^T (\text{approx } 0.1) \\ &= (\text{approx. } 0.00177) (W^{(3)})^T\end{aligned}$$

And it gets smaller and smaller as you go deeper.

When this happens, learning comes to a halt, especially in the deep layers. This is often called the

vanishing gradients problem.

Another way to look at it. Consider this simple, but deep network.



Start with the loss on the output side: $E(y, t)$

The gradient w.r.t. the input current of the output node is

$$\frac{\partial E}{\partial z^{(5)}} = y - t$$

Then, using backprop, we can compute a single formula for

$$\frac{\partial E}{\partial z^{(n)}} = (y - t) w^{(n)} \sigma'(z^{(n)})$$

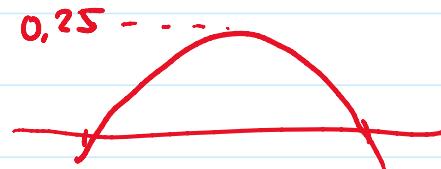
Going deeper...

$$\frac{\partial E}{\partial z^{(n)}} = (y - t) w^{(n)} \sigma'(z^{(n)}) w^{(n-1)} \sigma'(z^{(n-1)}) w^{(n-2)} \sigma'(z^{(n-2)}) \dots w^{(1)} \sigma'(z^{(1)})$$

What is the steepest slope that $\sigma(z)$ attains?

$$\sigma(z) = \sigma(z)(1 - \sigma(z))$$

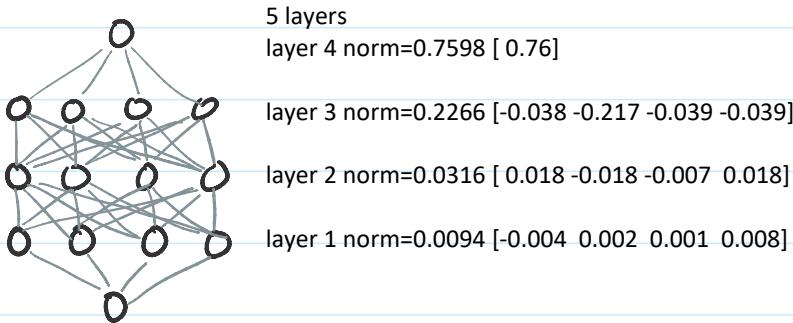
$$\text{for } 0 \leq \sigma(z) \leq 1$$



All else being equal, the gradient goes down by a factor of at least 4 each layer.

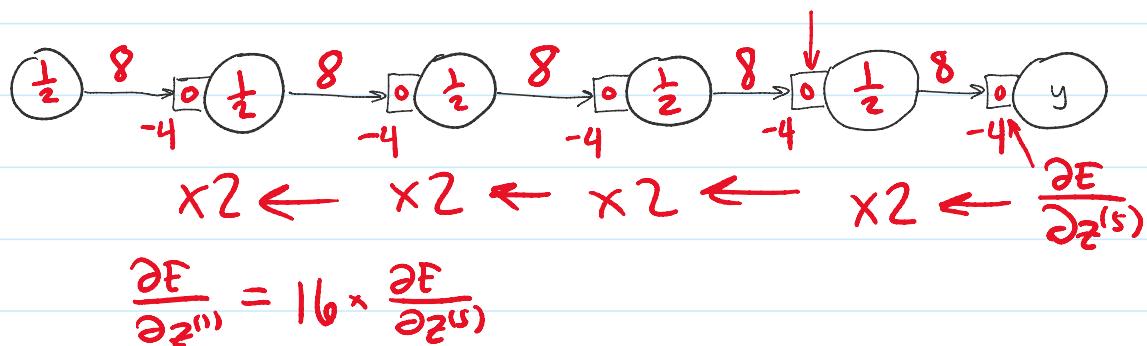
We can see this if we look at the norm of the gradients at each layer.

$$\text{L2. } \left\| \frac{\partial E}{\partial z^{(i)}} \right\|^2 = \sum_j \left(\frac{\partial E}{\partial z_j^{(i)}} \right)^2$$



Exploding Gradients

A similar, though less frequent phenomenon can result in very large gradients.



This situation is more rare since it only occurs when the weights are high and the biases compensate so that the input current lands in the sweet spot of the logistic curve.

