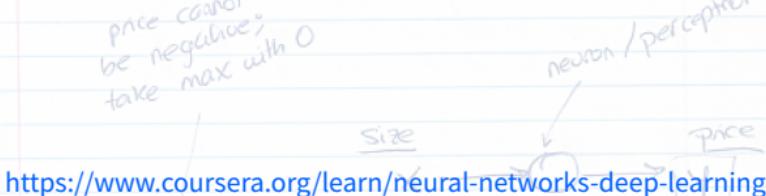


House Pricing

Course Notes for Neural Networks and Deep Learning



<https://www.coursera.org/learn/neural-networks-deep-learning>

$$y = \max(0, w_0 + w_1 x)$$

the REW part

weight

input

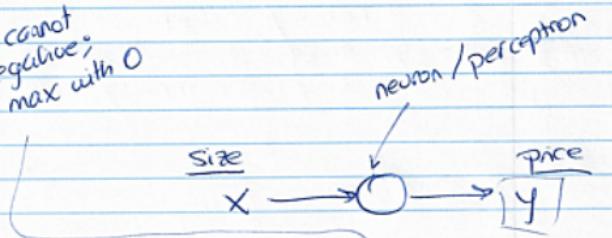
bias

Note: I am in no way affiliated to Coursera or the course authors.

House Pricing



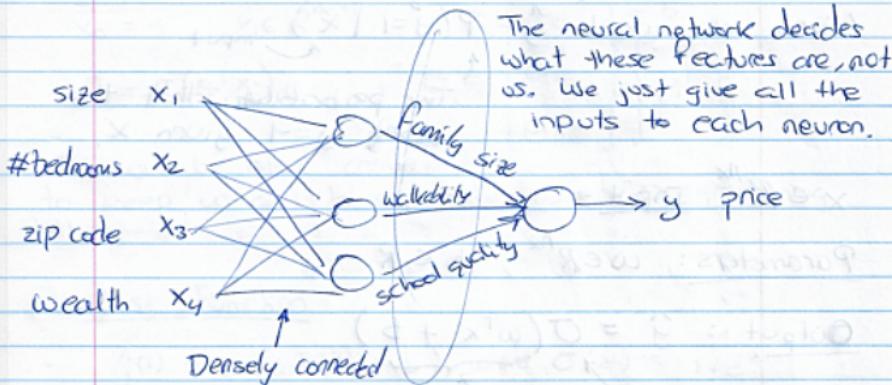
price cannot
be negative;
take max with 0



$$y = \max(0, wx + b)$$

The diagram shows the equation for the output y of the neuron. The term $wx + b$ is labeled "the REW part". The term $\max(0, \cdot)$ is labeled "bias". Arrows point from "input" to $wx + b$ and from "bias" to $\max(0, \cdot)$.

Multiple Variables



Notation

(x, y) : input/output pair $x \in \mathbb{R}^n$, $y \in \{0, 1\}$
(labeled input)

m training examples : $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

$$X = \left[\begin{array}{cccc} | & | & | & | \\ x^{(1)} & x^{(2)} & \cdots & x^{(m)} \\ | & | & | & | \end{array} \right] \underbrace{\}_{m^n}}$$

$$X \in \mathbb{R}^{n_x \times m}$$

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \cdots & y^{(m)} \end{bmatrix} \underbrace{\}_{1 \times m}}$$

$$Y \in \mathbb{R}^{1 \times m}$$

Logistic Regression

binary labels, $y \in \{0, 1\}$

Given x , want $\hat{y} = P(y=1 | x)$ input

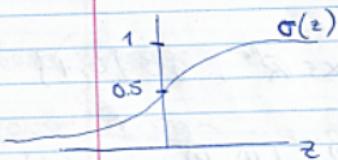
The probability that the label $y=1$ given x .

$x \in \mathbb{R}^n$ Input

Parameters: $w \in \mathbb{R}^n$, $b \in \mathbb{R}$

Output: $\hat{y} = \sigma(w^T x + b)$

Sigmoid function



$$\sigma(z) = \frac{1}{1+e^{-z}}$$

linear regression.
Will not produce a value $\hat{y} \in [0, 1]$, so unsuitable for probabilities.

$$\hookrightarrow \text{Large } z \rightarrow \sigma(z) \approx \frac{1}{1+0} \approx 1$$

$$\hookrightarrow \text{Large negative } z \rightarrow \sigma(z) \approx \frac{1}{1+\text{Big}} \approx 0$$

Maps the linear regression to a value $\hat{y} \in [0, 1]$ that can be interpreted as a probability.

Goal is to learn parameters w, b so that \hat{y} becomes a good estimate of the probability of $y=1$ given an input x .

principal or canonical test with binary classification

Alternative notation that merges w and b:

$$x_0 = 1, \quad x \in \mathbb{R}^{n+1}$$

$$\hat{y} = \sigma(w^T x)$$

$$w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_{n+1} \end{bmatrix} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_{n+1} \end{bmatrix} = b$$

It will be more convenient to keep w and b separate when programming this.

Cost Function

$$\hat{y}^{(i)} = \sigma(w^T x^{(i)} + b), \text{ where } \sigma(z) = \frac{1}{1+e^{-z}}$$

Given $\{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$
want $\hat{y}^{(i)} \approx y^{(i)}$ $x^{(i)}$ } i-th example
 $y^{(i)}$ }

Loss (error) function: → Measures how good our prediction (\hat{y}) is.

$$\text{One idea: } L(y, \hat{y}) = \frac{1}{2}(\hat{y} - y)^2$$

This is not very good in practice where we end up with a non-convex problem & the loss function we are trying to optimise has various local optima.



Gradient descent may therefore not find the global optima.

A better loss function that yields a convex optimisation problem:

$$L(\hat{y}, y) = -\left(y \log \hat{y} + (1-y) \log (1-\hat{y})\right)$$



global optima

Intuition behind the loss function in logistic regression:

$$L(\hat{y}, y) = - (y \log \hat{y} + (1-y) \log (1-\hat{y}))$$

If $y=1$:

$$\begin{aligned} L(\hat{y}, 1) &= - (\hat{y} \log \hat{y} + (1-\hat{y}) \log (1-\hat{y})) \\ &= - \log \hat{y} \end{aligned}$$

Since we are minimising $L(\hat{y}, y)$,
we want $-\log \hat{y}$ to be a large negative number
 \Leftrightarrow we want \hat{y} to be large
 $\Leftrightarrow \hat{y} = 1$, since sigmoid gives a value in $[0,1]$

If $y=0$:

$$L(\hat{y}, 0) = - \log (1-\hat{y})$$

Since we are minimising $L(\hat{y}, y)$,
we want $-\log (1-\hat{y})$ to be a large negative number
 \Leftrightarrow we want \hat{y} to be as small as possible
 $\Leftrightarrow \hat{y}=0$, since sigmoid gives a value in $[0,1]$.

Cost Function

The loss function is defined w.r.t a single training example.
To measure how well we are doing on the entire training set, we define the cost function.

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

| m training samples

cost function

$\underbrace{\hspace{10em}}$

The parameters we have learned

The average loss across the whole training set.

$$\frac{\partial}{\partial w} \frac{\partial}{\partial b}$$

Gradient Descent

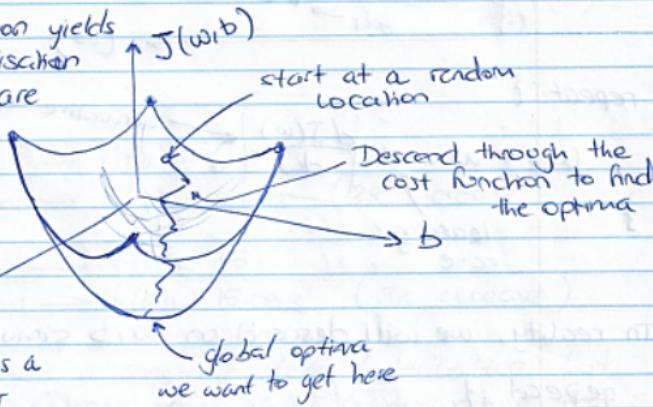
$$y = \sigma(w^T x + b), \quad \sigma(z) = \frac{1}{1 + e^{-z}}$$

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(y^{(i)}, \hat{y}^{(i)})$$

We want to find w, b that minimise $J(w, b)$

Our cost function yields a convex optimisation problem, so we are looking at a bowl with a single global optimum.

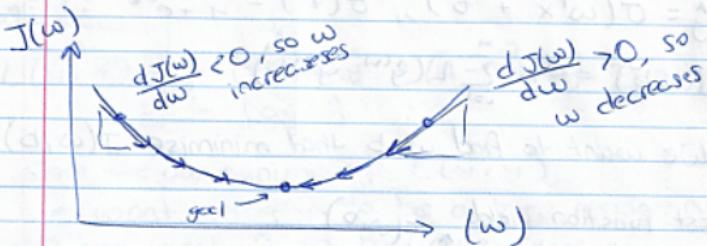
(Assuming this is a scalar just for illustrative purposes)



To start the optimisation procedure, w, b are initialised to either 0 or random values. 0 is the more popular choice for logistic regression.

If we used a different cost function such as the squared difference, the resulting optimisation problem would not be convex (have multiple local optima), and therefore gradient descent would not work really well.

scalar
Gradient descent on w (for illustration purposes):



repeat {

$$w := w - \alpha$$

$$\boxed{\frac{dJ(w)}{dw}} \leftarrow \text{derivative wrt } w$$

} learning rate

In reality, we will descend on w, b simultaneously:

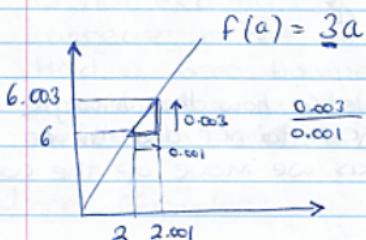
repeat {

$$w := w - \alpha \frac{\partial J(w, b)}{\partial w}$$

$$b := b - \alpha \frac{\partial J(w, b)}{\partial b}$$

}

Derivatives



$$f(a) = 3a$$

$$\frac{0.003}{0.001} = \frac{\text{height}}{\text{width}}$$

When we increase a by a small amount, $f(a)$ increases by $3x$ that amount.

$$\frac{dF(a)}{da} = \frac{df(a)}{da} = 3$$

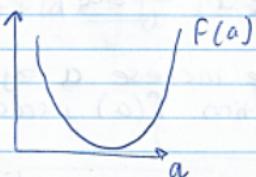
$$a = 2 \rightarrow f(a) = 6$$

$$a = 2.001 \rightarrow f(a) = 6.003 \quad (3x \text{ increase})$$

$$a = 5 \rightarrow f(a) = 15$$

$$a = 5.001 \rightarrow f(a) = 15.003 \quad (3x \text{ increase})$$

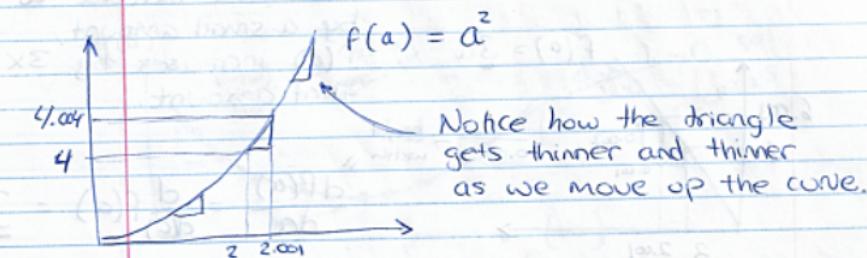
A line has a constant slope / derivative, but other functions may have a derivative that changes at every point. So the definition of the derivative must involve a very small increase.



$$\frac{dF(a)}{da} = \lim_{h \rightarrow 0} \frac{F(a+h) - F(a)}{h}$$

This is just the definition of slope (rise over run), but with a run that gets infinitesimally smaller.

A function with a varying derivative



$$\text{at } a = 2 \rightarrow f(a) = 4$$
$$a = 2.001 \rightarrow f(a) = 4.004001 \quad (4x \text{ increase})$$

$$a = 5 \rightarrow f(a) = 25$$
$$a = 5.001 \rightarrow f(a) = 25.010 \quad (10x \text{ increase})$$

The increase varies with the value of a .

$$\text{At } a = [2], \frac{d}{da} f(a) = [4] \quad (2 \cdot 2 \times)$$

$$\text{At } a = [5], \frac{d}{da} f(a) = [10] \quad (2 \cdot 5 \times)$$

$$\text{In general: } \frac{d}{da} f(a) = \frac{d}{da} a^2 = [2a]$$

What this means is that if we increase a by a small value ϵ , then the function $f(a)$ increases by:

$$a \rightarrow [a + \boxed{\epsilon}]$$
$$f(a) \rightarrow \underbrace{F(a)}_{\substack{\text{starting point}}} + \boxed{2a(\epsilon)}$$

↑ ↑
 increase

Note that when we bumped a from $a=2$, to $a=4$, $F(a)$ increased by 0.004001 , which is not exactly $2a = 4$. This is because the increase is not an infinitesimally small value. Had a been bumped by an infinitesimally small value, we would see an increase in $F(a)$ of exactly $2(a)$.

More Examples

$$f(a) = a^3$$

$$\frac{d}{da} f(a) = \boxed{3a^2}$$

$$\begin{aligned} a &= 2 \longrightarrow f(a) = 8 \\ +0.001 &\downarrow \\ a &= 2.001 \longrightarrow f(a) = 8.012 \\ 0.12 &= [0.001] \cdot (3 \cdot 2^2) \end{aligned}$$

$$\begin{aligned} F(a) &= \log_e(a) \\ &= \ln(a) \end{aligned}$$

$$\frac{d}{da} p(a) = \boxed{\frac{1}{a}}$$

$$\begin{aligned} a &= 2 \longrightarrow f(a) \approx 0.69315 \\ +0.0005 &\downarrow \\ a &= 2.001 \longrightarrow f(a) \approx 0.69365 \\ 0.0005 &= [0.001] \cdot \boxed{\frac{1}{2}} \end{aligned}$$

$$f(a)$$

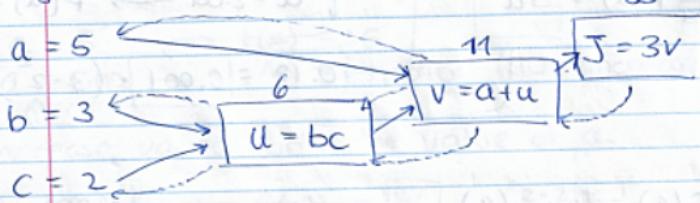
Computation Graph

A representation of a computation that breaks down smaller computations into nodes and where edges represent dependencies / ordering.

Example : $J(a,b,c) = 3(a+bc)$

1. Compute $U = bc$
2. Compute $V = a + U$
3. Compute $J = 3V$

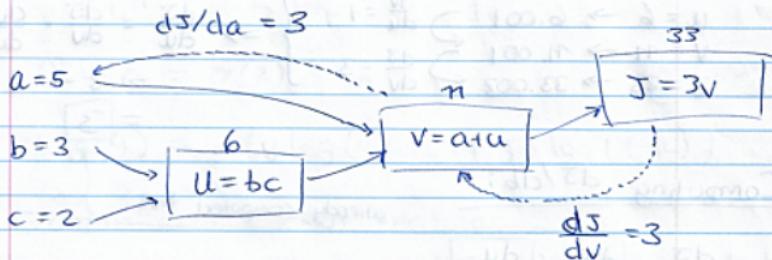
Visually :



The computation graph becomes useful when there is a value such as J that we want to optimise.

The output is easily computed with a forward or left-to-right pass. Derivatives can be computed with a backwards or right-to-left pass.

Derivatives with a Computation Graph



Computing dJ/dv :

$$J = 3v$$

$$\left. \begin{array}{l} v = 11 \rightarrow 11.001 \\ J = 33 \rightarrow 33.003 \end{array} \right\} \Rightarrow \frac{dJ}{dv} = 3$$

increase
 v by 0.001

$\xrightarrow{\quad}$
J increases by
 3×0.001

Computing dJ/da :

$$\left. \begin{array}{l} a = 5 \rightarrow 5.001 \\ v = 11 \rightarrow 11.001 \\ J = 33 \rightarrow 33.003 \end{array} \right\} \begin{array}{l} \frac{dv}{da} = 1 \\ \frac{dJ}{dv} = 3 \end{array} \Rightarrow \frac{dJ}{da} = \underline{\underline{3}}$$

Chain Rule:

$$\frac{dJ}{da} = \frac{dJ}{dv} \cdot \frac{dv}{da}$$

$$\underline{\underline{3}} = \underline{\underline{3}} \cdot \underline{\underline{1}}$$

Computing $\frac{dJ}{du}$:

$$\left. \begin{array}{l} u = 6 \rightarrow 6.001 \quad \frac{du}{du} = 1 \\ v = 11 \rightarrow 11.001 \quad \frac{dv}{du} = 3 \\ J = 33 \rightarrow 33.003 \quad \frac{dJ}{dv} = 3 \end{array} \right\} \Rightarrow \frac{dJ}{du} = \frac{dJ}{dv} \cdot \frac{dv}{du} = 3 \cdot 1 = \boxed{3}$$

Computing $\frac{dJ}{db}$:

already computed this

$$\frac{dJ}{db} = \frac{dJ}{du} \cdot \frac{du}{db} \quad b = 3 \rightarrow 3.001 \quad \frac{du}{db} = 2 \\ = 3 \cdot 2 \\ = 6$$

We can verify the result:

$$\left. \begin{array}{l} b = 3 \rightarrow 3.001 \quad \frac{du}{db} = 2 \\ u = 6 \rightarrow 6.002 \quad \frac{dv}{du} = 1 \\ v = 11 \rightarrow 11.002 \quad \frac{dJ}{dv} = 3 \\ J = 33 \rightarrow 33.006 \quad \frac{dJ}{du} = 3 \end{array} \right\} \frac{dJ}{db} = 2 \cdot 1 \cdot 3 = 6$$

$$\frac{dJ}{db} = \frac{dJ}{dv} \cdot \frac{dv}{du} \cdot \frac{du}{db} \\ = 3 \cdot 1 \cdot 2 \\ = 6$$

Computing $\frac{dJ}{dc}$:

$$\frac{dJ}{dc} = \frac{dJ}{du} \cdot \frac{du}{dc} = 3 \cdot 3 = 9$$

Logistic Regression Gradient Descent

$$z = w^T x + b$$

$$\hat{y} = a = \sigma(z)$$

$$L(a, y) = - (y \log(a) + (1-y) \log(1-a))$$

Example with 2 input features

$$\begin{array}{c}
 x_1 \\
 w_1 \\
 x_2 \\
 w_2 \\
 b
 \end{array}
 \rightarrow
 \boxed{z = w_1 x_1 + w_2 x_2 + b} \rightarrow
 \boxed{a = \sigma(z)} \rightarrow
 \boxed{L(a, y)}$$

$$\frac{\partial L(a, y)}{\partial z} = \frac{\partial L(a, y)}{\partial a} \cdot \frac{\partial a}{\partial z} = \frac{\partial L(a, y)}{\partial a} = -\frac{y}{a} + \frac{1-y}{1-a}$$

$$\frac{\partial L}{\partial w_1} = x_1 \cdot \frac{\partial L}{\partial z} = \left(-\frac{y}{a} + \frac{1-y}{1-a} \right) \cdot a(1-a)$$

$$\frac{\partial L}{\partial w_2} = x_2 \cdot \frac{\partial L}{\partial z} = a - y$$

x_1, x_2 are the input features.
 we want to tune the parameters w_1, w_2, b
 to minimise the loss function L .

Proceed backwards
 in the computation
 graph from the
 output to the input.

1. Compute dl/da
2. Compute dl/dz
3. Compute $dl/dw_1, dl/dw_2, dl/db$
4. Update:

$$w_1 := w_1 - \alpha \frac{dl}{dw_1}$$

$$w_2 := w_2 - \alpha \frac{dl}{dw_2}$$

$$b := b - \alpha \frac{dl}{db}$$

α : learning rate

Gradient Descent on m Examples

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(a^{(i)}, y^{(i)})$$

$$a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(w^T x^{(i)} + b)$$

Previously, we had seen how to compute $\frac{\partial w_1^{(i)}}{\partial w_1}$, $\frac{\partial w_2^{(i)}}{\partial w_1}$, $\frac{\partial b^{(i)}}{\partial w_1}$ for a single training example $(x^{(i)}, y^{(i)})$.

It turns out that to compute the derivative of the cost function wrt one of the parameters, we can simply average the derivative of the loss function wrt the same parameter over the m training examples. For example :

$$\frac{\partial}{\partial w_1} J(w, b) = \frac{1}{m} \underbrace{\sum_{i=1}^m \frac{\partial}{\partial w_1} L(a^{(i)}, y^{(i)})}_{\text{Average on all } m \text{ training examples}} \quad \begin{array}{l} \text{cost} \\ \text{loss} \end{array}$$

Already saw how to compute this term
 $\frac{\partial L(a^{(i)}, y^{(i)})}{\partial w_1}$

Example code is on the next page.

In practice, to speed up computations, we would use vectorisation instead of using explicit loops.

One single step of gradient descent over m training examples:

$$J = 0 \ ; \ dw_1 = 0 \ ; \ dw_2 = 0 \ ; \ db = 0$$

for $i = 1 \dots m$:

$$z^{(i)} = w^T x^{(i)} + b$$

$$a^{(i)} = \sigma(z^{(i)})$$

$$J += -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log (1-a^{(i)})]$$

$$dz^{(i)} = a^{(i)} - y^{(i)}$$

$$dw_1 += x_1^{(i)} dz^{(i)}$$

$$dw_2 += x_2^{(i)} dz^{(i)}$$

$$dw_n += x_n^{(i)} dz^{(i)}$$

$$db += dz^{(i)}$$

} For all n features.
In the examples so far
we used n=2.

$$J /= m$$

$$dw_1 //= m$$

$$dw_2 //= m$$

⋮

$$dw_n //= m$$

$$db //= m$$

} Dividing everything by m to
average out over the m
training examples.

// update!

$$w_1 := w_1 - \alpha dw_1$$

$$w_2 := w_2 - \alpha dw_2$$

⋮

$$w_n := w_n - \alpha dw_n$$

$$b := b - \alpha db$$

Vectorisation

The point of vectorisation is to remove explicit 'for' loops and instead process several elements (a vector) at a time using special CPU/GPU instructions.

Derivatives with Vectorisation during Gradient Descent

Instead of :

$$dw_1 = 0, dw_2 = 0, \dots, dw_{n_x} = 0$$

for $i = 1..M$:

$$\begin{aligned} dw_1 &= x_1^{(i)} dz^{(i)} \\ dw_2 &= x_2^{(i)} dz^{(i)} \end{aligned}$$

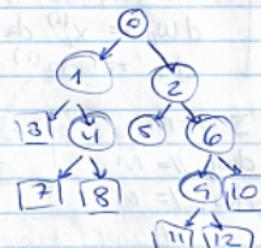
$$dw_{n_x} = x_{n_x}^{(i)} dz^{(i)}$$

$$dw_1 / = m$$

$$dw_2 / = m$$

$$\vdots$$

$$dw_{n_x} / = m$$



$$dw = \begin{bmatrix} dw_1 \\ dw_2 \\ \vdots \\ dw_{n_x} \end{bmatrix}$$

Auto expansion

$$dw = np. zeros (n_x, 1)$$

for loop:

for $i = 1..M$:

$$\begin{aligned} dw &+= X^{(i)} dz^{(i)} \\ dw / &= m \end{aligned}$$

$$\left\{ \begin{bmatrix} dw_1 \\ dw_2 \\ \vdots \\ dw_{n_x} \end{bmatrix} \right. \left. + \begin{bmatrix} x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_{n_x}^{(i)} \end{bmatrix} \begin{bmatrix} dz^{(i)}_1 \\ dz^{(i)}_2 \\ \vdots \\ dz^{(i)}_{n_x} \end{bmatrix} \right\}$$

These are scalars, but are automatically expanded

Vectorizing Logistic Regression

To perform a forward pass of logistic regression given m training examples, we need to compute:

$$\begin{aligned} z^{(1)} &= \mathbf{w}^T \mathbf{x}^{(1)} + b & z^{(2)} &= \mathbf{w}^T \mathbf{x}^{(2)} + b & \dots & z^{(m)} = \mathbf{w}^T \mathbf{x}^{(m)} + b \\ a^{(1)} &= \sigma(z^{(1)}) & a^{(2)} &= \sigma(z^{(2)}) & \dots & a^{(m)} = \sigma(z^{(m)}) \end{aligned}$$

Instead of iterating over the m training examples explicitly with a for loop, we define a huge matrix:

$$\mathbf{X} = \left[\begin{array}{c|c|c|c} & | & | & | \\ \mathbf{x}^{(1)} & \mathbf{x}^{(2)} & \dots & \mathbf{x}^{(m)} \\ | & | & | & | \end{array} \right] \quad \mathbf{X} \in \mathbb{R}^{n_x \times m}$$

Then:

$$\mathbf{z} = \begin{bmatrix} z^{(1)} & z^{(2)} & \dots & z^{(m)} \end{bmatrix} = \mathbf{w}^T \mathbf{X} + \begin{bmatrix} b & b & \dots & b \end{bmatrix}$$

$1 \times m \qquad \qquad \qquad 1 \times n_x \qquad n_x \times m \qquad 1 \times m$

Where \mathbf{w}^T is the row vector $\mathbf{w}^T = [\omega^{(1)} \ \omega^{(2)} \ \dots \ \omega^{(n_x)}]$

$$\begin{bmatrix} z^{(1)} & z^{(2)} & \dots & z^{(m)} \end{bmatrix} = \begin{bmatrix} \omega^{(1)} & \omega^{(2)} & \dots & \omega^{(n_x)} \end{bmatrix} \begin{bmatrix} | & | & | \\ \mathbf{x}^{(1)} & \mathbf{x}^{(2)} & \dots & \mathbf{x}^{(m)} \\ | & | & | & | \end{bmatrix} + \begin{bmatrix} b & b & \dots & b \end{bmatrix}$$

In code:

$$\underline{\mathbf{z}} = \text{np. dot}(\mathbf{w}, \mathbf{T}, \mathbf{X}) + b \quad \text{scalar } b, \text{ automatically broadcasted}$$

Finally, we would compute:

$$\mathbf{A} = \begin{bmatrix} a^{(1)} & a^{(2)} & \dots & a^{(m)} \end{bmatrix} = \sigma(\mathbf{z})$$

Vectorising the Gradient Computation

1. Instead of computing $dz^{(i)}$ over the m training examples with a loop:

$$dz^{(1)} = a^{(1)} - y^{(1)}, dz^{(2)} = a^{(2)} - y^{(2)}, \dots, dz^{(m)} = a^{(m)} - y^{(m)}$$

define:

$$dZ = A - Y$$

$$[dz^{(1)} \ dz^{(2)} \ \dots \ dz^{(m)}] = [a^{(1)} - y^{(1)}, a^{(2)} - y^{(2)}, \dots, a^{(m)} - y^{(m)}]$$

2. To compute db ,

Instead of:

$$\begin{aligned} db &= 0 \\ db + &= dz^{(1)} \\ db + &= dz^{(2)} \\ &\vdots \\ db + &= dz^{(m)} \\ db / &= m \end{aligned}$$

Do this:

$$\begin{aligned} db &= \frac{1}{m} \sum_{i=1}^m dz^{(i)} \\ &= \frac{1}{m} \text{np.sum}(dz) \end{aligned}$$

3. To compute dW ,

Instead of: vectors scalars

$$\begin{aligned} dW &= 0 \\ dW + &= x^{(1)} dz^{(1)} \\ dW + &= x^{(2)} dz^{(2)} \\ &\vdots \\ dW + &= x^{(m)} dz^{(m)} \\ dW / &= m \end{aligned}$$

Do this:

$$dW = \frac{1}{m} X dZ^T$$

To see why the vectorised computation of $d\mathbf{w}$ works:

$$d\mathbf{w} = \frac{1}{m} \times d\mathbf{z}^T$$

$$= \frac{1}{m} \begin{bmatrix} |^{(1)} & |^{(2)} & |^{(m)} \\ X^{(1)} & X^{(2)} & \cdots & X^{(m)} \\ | & | & | & | \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ dz^{(2)} \\ \vdots \\ dz^{(m)} \end{bmatrix}$$

$$= \frac{1}{m} \underbrace{\left[X^{(1)} dz^{(1)} + X^{(2)} dz^{(2)} + \cdots + X^{(m)} dz^{(m)} \right]}_{n \times 1}$$

Example:

$$\begin{bmatrix} X_1^{(1)} & X_1^{(2)} \\ X_2^{(1)} & X_2^{(2)} \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ dz^{(2)} \end{bmatrix} = \begin{bmatrix} X_1^{(1)} dz^{(1)} + X_1^{(2)} dz^{(2)} \\ X_2^{(1)} dz^{(1)} + X_2^{(2)} dz^{(2)} \end{bmatrix}$$

$$= \begin{bmatrix} X_1^{(1)} dz^{(1)} \\ X_2^{(1)} dz^{(1)} \end{bmatrix} + \begin{bmatrix} X_1^{(2)} dz^{(2)} \\ X_2^{(2)} dz^{(2)} \end{bmatrix} = dz^{(1)} \begin{bmatrix} X_1^{(1)} \\ X_2^{(1)} \end{bmatrix} + dz^{(2)} \begin{bmatrix} X_1^{(2)} \\ X_2^{(2)} \end{bmatrix}$$

$$= \underbrace{\frac{X^{(1)}}{\text{vector}}}_{\text{scalar}} \underbrace{\frac{dz^{(1)}}{\text{vector}}}_{\text{scalar}} + \underbrace{\frac{X^{(2)}}{\text{vector}}}_{\text{scalar}} \underbrace{\frac{dz^{(2)}}{\text{vector}}}_{\text{scalar}}$$

Putting it all together to implement logistic regression:

$$z = w^T X + b = \text{np.dot}(w, T(X)) + b$$

$$A = \sigma(z)$$

$$dz = A - Y$$

$$dw = \frac{1}{m} \times dz^T$$

$$db = \frac{1}{m} \text{np.sum}(dz)$$

$$w = w - \alpha dw$$

$$b = b - \alpha db$$

This implements one step of forward propagation and backward propagation on all m training samples.

Note that we will still need a for loop over the number of training iterations.

Logistic Regression Cost Function

$$g = \sigma(w^T x + b), \quad \sigma(z) = \frac{1}{1+e^{-z}}$$

Interpret $\hat{y} = p(y=1|x)$, which in other words is:

$$\left. \begin{array}{l} \text{if } y=1: \quad p(y|x) = \hat{y} \\ \text{if } y=0: \quad p(y|x) = 1-\hat{y} \end{array} \right\} p(y|x)$$

Next, we summarise these two equations in a single equation:

$$p(y|x) = \hat{y}^y (1-\hat{y})^{(1-y)}$$

$$\text{If } y=1: p(y|x) = \hat{y}^1 (1-\hat{y})^0 = \hat{y}$$

$$\text{If } y=0: p(y|x) = \hat{y}^0 (1-\hat{y})^{(1-0)} = 1-\hat{y}$$

Finally, we are going to maximise $\log p(y|x)$. Note that since the log function is monotonically increasing, maximising $\log p(y|x)$ is the same as maximising $p(y|x)$.

$\log p(y|x)$

$$\log p(y|x) = \log \hat{y}^y (1-\hat{y})^{(1-y)}$$

Maximise
 $\log p(y|x)$
= minimise
negative
loss

$$= y \log \hat{y} + (1-y) \log (1-\hat{y})$$

This defines the loss function for a single example. Now, let's define the cost function for M examples.

Assuming that the training examples are drawn independently and are identically distributed:

$$p(\text{labels in training set}) = \prod_{i=1}^M p(y^{(i)} | x^{(i)})$$

To carry out maximum likelihood estimation, we find the parameters that maximise the probability of the observations in the training set. Again, instead of maximising the probability directly, we maximise its log:

$$\begin{aligned}\log P(\text{labels}) &= \sum_{i=1}^m \underbrace{\log P(y^{(i)}, x^{(i)})}_{-L(g^{(i)}, y^{(i)})} \\ &= -\sum_{i=1}^m L(g^{(i)}, y^{(i)})\end{aligned}$$

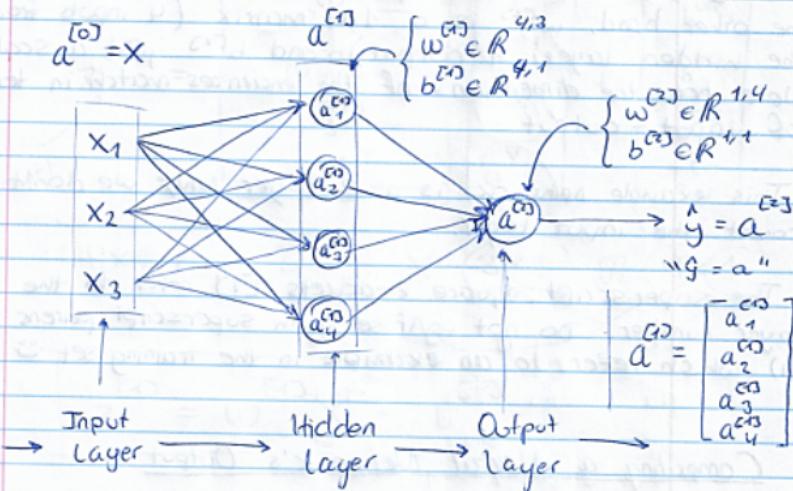
Cost:
(minimise) $J(w, b) = \frac{1}{m} \sum_{i=1}^m L(g^{(i)}, y^{(i)})$

Minimising the cost is maximising the likelihood, so the negative sign $-\sum$ disappears.

The $\frac{1}{m}$ is added for convenience, so that quantities are better scaled. The cost function ends up being the average of the loss function across the M training examples.

SHALLOW Neural Networks I

Neural Network Representation



"Hidden" refers to the fact that the values for the hidden layer are not observed in the training set. The training set only contains (input, output) pairs.

" a " stands for activations. They are the output of each layer. The inputs can be thought as the first set of activations.

Note the dimension of things. In logistic regression, " a " is a scalar. Here, $a^{(0)}$ is a 3D vector, $a^{(1)}$ is a 4D vector, and $a^{(2)}$ (the final output) is a scalar.

Similarly, every non-input layer is associated with a set of parameters W, b . W need no longer be a vector, nor " b " a scalar. In the example, $W^{[1]}$ is a 4×3 matrix (the layer has 3 inputs & 4 outputs), and b is a 4×1 matrix (a 4D column vector). On the other hand, $W^{[2]}$ is a 1×4 matrix (4 inputs from the hidden layer, 1 output) and $b^{[2]}$ just a scalar. Note how the dimensions of the matrices match in terms of input-output.

This example network is a 2 layer NN; we don't count the input layer.

The superscript square brackets [1] refer to the layer number. Do not confuse with superscript parens (1) which refer to an example in the training set \tilde{x}

Computing a Neural Network's Output

The neural network is basically going to perform one logistic regression per node / neuron.

Hidden Layer Computations

$$\begin{aligned} z_1^{[1]} &= W_1^{[1]T} x + b_1^{[1]}, & a_1^{[1]} &= \sigma(z_1^{[1]}) \\ z_2^{[1]} &= W_2^{[1]T} x + b_2^{[1]}, & a_2^{[1]} &= \sigma(z_2^{[1]}) \\ z_3^{[1]} &= W_3^{[1]T} x + b_3^{[1]}, & a_3^{[1]} &= \sigma(z_3^{[1]}) \\ z_4^{[1]} &= W_4^{[1]T} x + b_4^{[1]}, & a_4^{[1]} &= \sigma(z_4^{[1]}) \end{aligned}$$

The notation is:

- $[L] \leftarrow$ layer
- $z_i \leftarrow$ node inside the layer

And of course, we are going to vectorise the whole computation. This is done by putting all the weight vectors into a matrix (this is done for each layer):

$$z^{[1]} = \underbrace{\begin{bmatrix} w_1^{[0]T} \\ w_2^{[0]T} \\ w_3^{[0]T} \\ w_4^{[0]T} \end{bmatrix}}_{4 \times 1} \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}}_{3 \times 1} + \underbrace{\begin{bmatrix} b_1^{[0]} \\ b_2^{[0]} \\ b_3^{[0]} \\ b_4^{[0]} \end{bmatrix}}_{4 \times 1}$$

No subscript!

\equiv

$4 \times 1 \quad 4 \times 3 \quad 3 \times 1 \quad 4 \times 1$

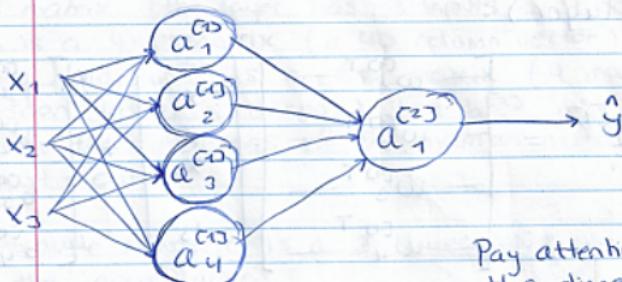
$$z^{[1]} = w^{[1]} x^{[0]} + b^{[1]}$$

Similarly for the activations:

$$a^{[1]} = \begin{bmatrix} a_1^{[0]} \\ a_2^{[0]} \\ a_3^{[0]} \\ a_4^{[0]} \end{bmatrix} = \begin{bmatrix} \sigma(z_1^{[0]}) \\ \sigma(z_2^{[0]}) \\ \sigma(z_3^{[0]}) \\ \sigma(z_4^{[0]}) \end{bmatrix}$$

$$a^{[1]} = \sigma(z^{[1]})$$

Finally, we would go about doing the same thing for the output layer. The computations are all summarised here:



Pay attention to
the dimension
of things!

Given input x :

$$z^{[1]} = W^{[1]} a^{[0]} + b^{[1]}$$

4×1 4×3 3×1 4×1

$$a^{[1]} = \sigma(z^{[1]})$$

4×1 4×1

$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$

1×1 1×4 4×1 1×1

$$a^{[2]} = \sigma(z^{[2]})$$

1×1 1×1

Vectorising Across Multiple Examples

Now we want to compute the output of the NN, but for all training examples:

$$\begin{aligned} x^{(1)} &\rightarrow g^{(1)} = a^{(2)}(1) \\ x^{(2)} &\rightarrow g^{(2)} = a^{(2)}(2) \\ &\vdots \\ x^{(m)} &\rightarrow g^{(m)} = a^{(2)}(m) \end{aligned}$$

(Without vectorisation, we would just repeat what we did in the previous section but for all examples:

for $i=1..M$:

$$\begin{aligned} z^{(1)}(i) &= w^{(1)} x^{(i)} + b^{(1)} \\ a^{(1)}(i) &= \sigma(z^{(1)}(i)) \end{aligned} \quad \left. \begin{array}{l} \text{first} \\ \text{layer} \end{array} \right\}$$

$$\begin{aligned} z^{(2)}(i) &= w^{(2)} a^{(1)}(i) + b^{(2)} \\ a^{(2)}(i) &= \sigma(z^{(2)}(i)) \end{aligned} \quad \left. \begin{array}{l} \text{second} \\ \text{layer} \end{array} \right\}$$

To vectorise this, we construct one giant matrix with all training examples stacked horizontally:

$$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \cdots x^{(m)} \\ | & | & | \end{bmatrix} \quad X \in \mathbb{R}^{(n_x, m)}$$

And then we compute:

matrices

$$\begin{aligned} z^{(1)} &= w^{(1)} X + b^{(1)} \\ A^{(1)} &= \sigma(z^{(1)}) \\ z^{(2)} &= w^{(2)} A^{(1)} + b^{(2)} \\ A^{(2)} &= \sigma(z^{(2)}) \end{aligned}$$

To visualise this:

$$z^{(l)} = w^{(l)} X + b^{(l)}$$

$$z^{(l)} = \begin{bmatrix} | & | & | & | \\ z^{(l)(1)} & z^{(l)(2)} & \dots & z^{(l)(m)} \\ | & | & | & | \end{bmatrix}$$

$$z^{(l)} = \begin{bmatrix} -w_1^{(l)T} - \\ -w_2^{(l)T} - \\ -w_n^{(l)T} - \end{bmatrix} \begin{bmatrix} | & | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | & | \end{bmatrix} + \begin{bmatrix} b_1^{(l)} \\ b_2^{(l)} \\ \vdots \\ b_m^{(l)} \end{bmatrix}$$

Here, ' n ' is the number of neurons in the layer, and ' m ' is the number of training examples.

Finally:

$$A^{(l)} = \sigma(z^{(l)}) = \begin{bmatrix} | & | & | & | \\ a^{(l)(1)} & a^{(l)(2)} & \dots & a^{(l)(m)} \\ | & | & | & | \end{bmatrix}$$

Interpretation:

activation
 $a_{i,j}$
ith neuron
of the layer

j-th training
example

$$A^{(l)} \in \mathbb{R}^{n \times m}$$

$$A^{(l)} = \begin{bmatrix} | & | & | \\ \text{---} & \text{---} & \text{---} \\ | & | & | \end{bmatrix}$$

neuron

example

More visualisation follows. Given that we have to vectorise:

$$z^{(1)}(1) = w^{(1)} x^{(1)} + b^{(1)}$$

$$z^{(1)}(2) = w^{(1)} x^{(2)} + b^{(1)}$$

$$z^{(1)}(3) = w^{(1)} x^{(3)} + b^{(1)}$$

Let's pretend that $b^{(1)} = 0$, to make things easier. Then:

$$w^{(1)} = \begin{bmatrix} \quad \\ \quad \\ \quad \end{bmatrix}$$

$$w^{(1)} x^{(1)} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$w^{(1)} x^{(2)} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$w^{(1)} x^{(3)} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

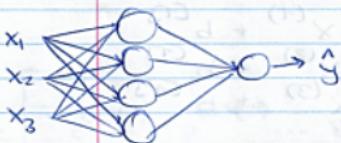
$$w^{(1)} X = w^{(1)} \underbrace{\begin{bmatrix} 1 & 1 & 1 \\ x^{(1)} & x^{(2)} & x^{(3)} \\ 1 & 1 & 1 \end{bmatrix}}_X = \underbrace{\begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}}_{z^{(1)}}$$

$$= \begin{bmatrix} z^{(1)(1)} & z^{(1)(2)} & z^{(1)(3)} \\ | & | & | \end{bmatrix} = z^{(1)}$$

$$\xrightarrow{\quad} \left\{ + b^{(1)} \quad + b^{(1)} \quad + b^{(1)} \right\}$$

With a non-zero $b^{(1)}$, $b^{(1)}$ is broadcasted across the columns of $z^{(1)}$.

To summarise:



$$X = \begin{bmatrix} | & | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | & | \end{bmatrix}$$

$$A^{(l)} = \begin{bmatrix} | & | & | & | \\ a^{(l)(1)} & a^{(l)(2)} & \dots & a^{(l)(m)} \\ | & | & | & | \end{bmatrix}$$

Without vectorisation:

for $i = 1..m$:

$$z^{(l)(i)} = w^{(l)} X^{(i)} + b^{(l)}$$

$$a^{(l)(i)} = \sigma(z^{(l)(i)})$$

$$z^{(l+1)(i)} = w^{(l+1)} A^{(l)(i)} + b^{(l+1)}$$

$$a^{(l+1)(i)} = \sigma(z^{(l+1)(i)})$$

With vectorisation:

$$\{ z^{(l)} = w^{(l)} X + b^{(l)} \}$$

$$\{ A^{(l)} = \sigma(z^{(l)}) \}$$

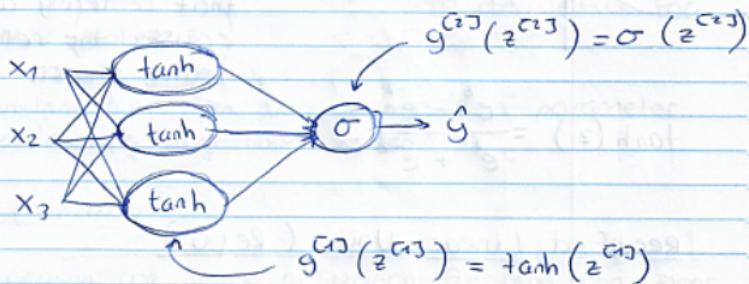
$$\{ z^{(l+1)} = w^{(l+1)} A^{(l)} + b^{(l+1)} \}$$

$$\{ A^{(l+1)} = \sigma(z^{(l+1)}) \}$$

Each of these four lines correspond to the lines above but vectorised across all m training examples.

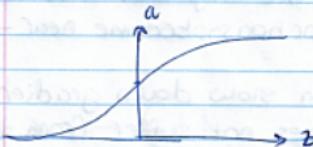
Activation Functions

One can use activation functions other than the sigmoid in the layers of a neural network. The activation function will be the same for all neurons in a given layer, but each layer can use a different one.



We use ' g ' to denote the activation function, and square superscripts to denote the layer.

Sigmoid

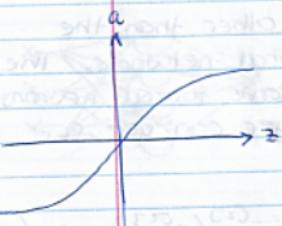


$$a = \frac{1}{1 + e^{-z}}$$

Use this only for the output layer. If your labels are $y \in \{0, 1\}$, then the sigmoid is convenient because it outputs $0 \leq y \leq 1$.

For hidden layers, the tanh function is a better choice.

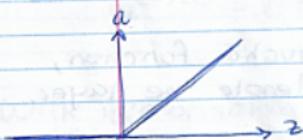
Hyperbolic Tangent



$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

Use this rather than the sigmoid for hidden layers. With tanh, the activations of a layer will have a mean closer to 0, and that centering around 0 causes the network to learn faster.

Rectified Linear Unit (Relu)



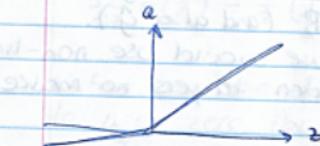
$$a = \max(0, z)$$

This is a good default choice.

When z is very large or very small, the derivatives of the sigmoid and tanh functions become near-zero.

These near-zero derivatives then slow down gradient descent. The relu function does not suffer from this problem, because the derivative is 0 for $z > 0$. Of course, for $z < 0$ the derivative of the relu is 0, but this does not usually cause problems in practice. If it does, try the leaky relu.

Leaky Relu



Also a good default choice, though apparently not as popular as the Relu.

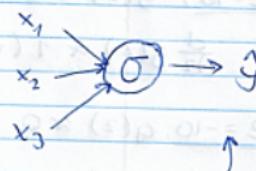
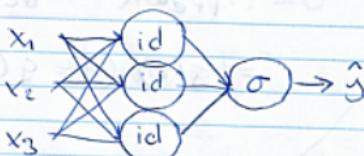
Fixes the problem with the zero derivative for $z < 0$.

The constant 0.01 can also be made a parameter of the learning algorithm instead.

Non-linearity

The reason we make activation functions non-linear is to make neural networks more powerful / expressive. If we were to use the identity as the activation function, then all we would be doing is computing a linear combination of the inputs, no matter how many hidden layers we add. This defeats deep networks.

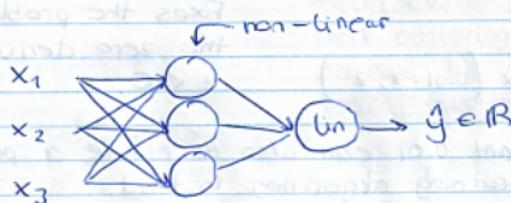
For instance:



This is not more expressive than the network on the right because the composition of two linear functions is itself a linear function.

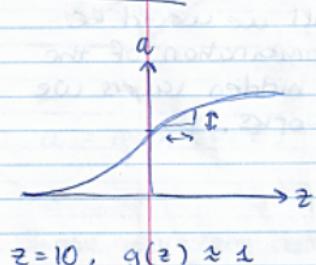
One case where we would use a linear activation function like the identity is when computing a regression, that is when $y \in \mathbb{R}$ (and also g).

Even in this case, however, we would use non-linear activation functions in the hidden layers to make the network more expressive.



Derivatives of Activation Functions

Sigmoid



$$z=10, g(z) \approx 1$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

$$\begin{aligned} \frac{d}{dz} g(z) &= \text{slope of } g(x) \text{ at } z \\ &= \frac{1}{1+e^z} \left(1 - \frac{1}{1+e^z} \right) \\ &= g(z) (1 - g(z)) \end{aligned}$$

$$\frac{d}{dz} g(z) \approx 1(1-1) \approx 0$$

$$z=-10, g(z) \approx 0$$

$$\frac{d}{dz} g(z) \approx 0(1-0) \approx 0$$

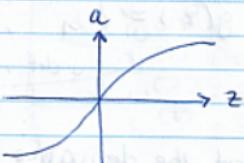
Collapse!

$$z=0, g(z)=\frac{1}{2}$$

$$\frac{d}{dz}g(z) = \frac{1}{2}(1-\frac{1}{2}) = \frac{1}{4}$$

Notice how the derivative of the sigmoid collapses to near zero for large positive/negative values of z . This is what we talked about when introducing the ReLU.

Tanh



$$g(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\begin{aligned}\frac{d}{dz}g(z) &= 1 - (\tanh(z))^2 \\ &= 1 - g^2(z)\end{aligned}$$

$$z=10, \tanh(z) \approx 1$$

$$\frac{d}{dz}g(z) \approx 0$$

Collapses for large $\pm z$ just like the sigmoid.

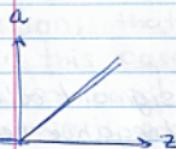
$$z=-10, \tanh(z) \approx -1$$

$$\frac{d}{dz}g(z) \approx 0$$

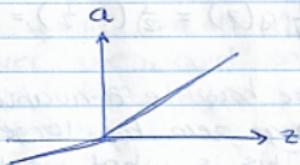
$$z=0, \tanh(z) \approx 0$$

$$\frac{d}{dz}g(z) = 1$$

ReLU and Leaky REU



$$g(z) = \max(0, z)$$



$$g(z) = \max(0.01z, z)$$

$$\frac{d}{dz} g(z) = \begin{cases} 0, & z < 0 \\ 1, & z \geq 0 \\ \text{undef}, & z=0 \end{cases}$$

$$\frac{d}{dz} g(z) = \begin{cases} 0.01, & z < 0 \\ 1, & z \geq 0 \\ \text{undef}, & z=0 \end{cases}$$

In software, we'll just pretend that the derivative is just 1 or 0 when $z=0$. This is technically a subgradient of the activation function, which works just fine for our purposes.

Gradient Descent for Neural Networks

Parameters : $w^{[1]} \in \mathbb{R}^{n^{[1]} \times n^{[2]}}$
 $b^{[1]} \in \mathbb{R}^{n^{[1]} \times 1}$
 $w^{[2]} \in \mathbb{R}^{n^{[2]} \times n^{[3]}}$
 $b^{[2]} \in \mathbb{R}^{n^{[2]} \times 1}$

Cost Function: $J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]})$
 $= \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i)$
 $\ell_a^{[2]}$

Gradient Descent:

repeat {

// Forward pass

compute predictions $\{\hat{y}^{(i)} \mid i=1..m\}$

// Backprop, compute derivatives

$$dw^{[2]} = \frac{dJ}{dw^{[2]}} , db^{[2]} = \frac{dJ}{db^{[2]}}$$

$$dw^{[2]} = \frac{dJ}{dw^{[2]}} , db^{[2]} = \frac{dJ}{db^{[2]}}$$

// Update parameters

$$w^{[1]} := w^{[1]} - \alpha dw^{[1]}$$

$$b^{[1]} := b^{[1]} - \alpha db^{[1]}$$

$$w^{[2]} := w^{[2]} - \alpha dw^{[2]}$$

$$b^{[2]} := b^{[2]} - \alpha db^{[2]}$$

Forward Propagation

$$A^{[1]} = g^{[1]} \quad (z^{[1]} = w^{[1]} X + b^{[1]})$$

$$A^{[2]} = g^{[2]} \quad (z^{[2]} = w^{[2]} A^{[1]} + b^{[2]})$$

Back propagation

$$dz^{[2]} = A^{[2]} = Y$$

$$dw^{[2]} = \frac{1}{m} dz^{[2]} A^{[1]T}$$

$$db^{[2]} = \frac{1}{m} np.sum(dz^{[2]}, axis=1, keepdims=True)$$

$$dz^{[1]} = w^{[2]T} dz^{[2]} \text{ /* } g^{[1]}(z^{[1]})$$

$$dw^{[1]} = \frac{1}{m} dz^{[1]} X^T$$

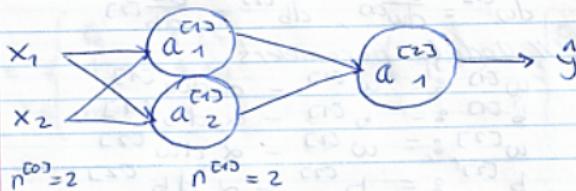
$$db^{[1]} = \frac{1}{m} np.sum(dz^{[1]}, axis=1, keepdims=True)$$

/* element-wise product

Random Initialization

When training a neural network, we must initialise the weights randomly instead of setting them to 0, otherwise the training will not work. This was OK for logistic regression, but not for networks with more than one neuron. The bias vectors can still be zero.

Consider:



Now initialise:

$$w^{(1)} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}, b^{(1)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$
$$w^{(2)} = \begin{bmatrix} 0 & 0 \end{bmatrix}$$

Now, for any given input, the two neurons in the hidden layer will compute the same function, that is:

$$a_1^{(1)} = a_2^{(1)}$$

Consequently, when we compute backpropagation, the derivatives will be the same:

$$dz_1^{(1)} = dz_2^{(1)}$$

It can then be shown by induction that after every single training iteration, the two hidden

units will still be computing the same function.

$$dw = \begin{bmatrix} u & v \\ u & v \end{bmatrix} \quad w^{[1]} := w^{[1]} - \alpha dw$$

$$w^{[1]} = \begin{bmatrix} \text{---} \\ \text{---} \end{bmatrix} \quad \left\{ \begin{array}{l} \text{identical} \\ \text{rows} \end{array} \right.$$

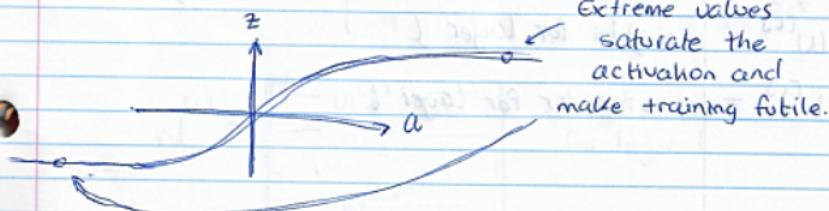
With a zero initialisation, all neurons in a layer are symmetric, and this symmetry never breaks throughout the training.

Instead, initialise weights randomly. Bias vectors can still be zero because they don't suffer from the symmetry problem:

$$\begin{aligned} w^{[1]} &= \text{np.random.randn}(2, 2) \cdot 0.01 \\ b^{[1]} &= \text{np.zeros}(2, 1) \\ w^{[2]} &= \dots \\ b^{[2]} &= \dots \end{aligned}$$

small value!

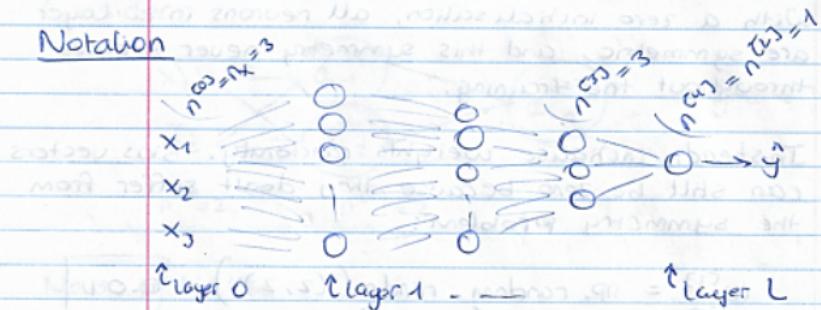
The constant 0.01 is to give initial weights small values. This is to prevent saturation of sigmoid and tanh functions. This matters even if you only use sigmoid and tanh in the output layer, and ReLUs in the hidden ones.



Finally, note that the constant 0.01 will likely work only for shallow networks. As you add more layers to the network, you will want to make this constant smaller.

Deep Neural Network

Notation



$$L = \# \text{ layers}$$

$$n^{[l]} = \# \text{ units in layer } l$$

$$a^{[l]} = \text{activations of layer } l \quad (a^{[l]} = g^{[l]}(z^{[l]}))$$

$$a^{[0]} = X = \text{input}$$

$$a^{[L]} = \hat{y} = \text{output}$$

$$w^{[l]} = \text{weights for layer } l$$

$$b^{[l]} = \text{bias vector for layer } l$$

Forward Propagation

There is nothing particularly new here; just more layers.

Given an input, we need to compute:

for $l = 1 \dots L$:

$$\begin{aligned} z^{[l]} &= w^{[l]} a^{[l-1]} + b^{[l]} \\ a^{[l]} &= g^{[l]}(z^{[l]}) \end{aligned}$$

Where:

$a^{[0]} = x = \text{input}$

$a^{[L]} = y = \text{output}$

$g^{[l]}$ = activation function of layer l

$w^{[l]}, b^{[l]}$ = weights and bias

We can then vectorise over all 'm' training examples like we had done before:

for $l = 1 \dots L$: ← Can't get rid of this loop

$$\begin{aligned} z^{[l]} &= w^{[l]} A^{[l-1]} + b^{[l]} \\ A^{[l]} &= g^{[l]}(z^{[l]}) \end{aligned}$$

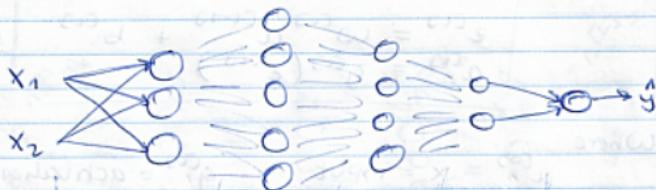
Where

$$X = A^{[0]} = \left[\begin{array}{c|c|c|c} 1 & x^{(1)} & x^{(2)} & \cdots & x^{(m)} \\ \hline 1 & | & | & \cdots & | \end{array} \right] \text{Training examples input.}$$

$$W^{[l]} = \left[\begin{array}{c|c} w_1^{[l]T} & \\ \hline w_2^{[l]T} & \\ \hline \vdots & \\ \hline w_n^{[l]T} & \end{array} \right] \text{Weight matrix for layer } l$$

The explicit for-loop on the layers is still necessary because we need the output of one layer to compute the output of the following layer. There does not appear to be a way to vectorise this.

Matrix Dimensions



$$n^{C1} = n_x = 2 \quad n^{C1} = 3 \quad n^{C1} = 5 \quad n^{C1} = 4 \quad n^{C1} = 2 \quad n^{C1} = 1$$

$$z^{C1} = w^{C1} x + b^{C1}$$

$$\begin{aligned} (3,1) &= (3,2)(2,1) + (3,1) \\ (n^{C1}, 1) &= (n^{C1}, n^{C1})(n^{C1}, 1) + (n^{C1}, 1) \end{aligned}$$

$\boxed{}$

$$w^{C1} = (n^{C1}, n^{C1})$$

$$b^{C1} = (n^{C1}, 1)$$

$$dw^{C1} = (n^{C1}, n^{C1})$$

$$db^{C1} = (n^{C1}, 1)$$

A matrix maps inputs to outputs. The number of columns is the number of inputs, and the number of rows is the number of outputs:

$$W^{[l]} \rightarrow \begin{pmatrix} n^{c_l}, n^{c_{l-1}} \end{pmatrix}$$

n^{c_l} outputs $n^{c_{l-1}}$ inputs

Vectorised Version

The vectorised version does not change much. Now we have:

$$\tilde{z}^{[l]} = \begin{bmatrix} z^{[l]}(1) & z^{[l]}(2) & \dots & z^{[l]}(m) \end{bmatrix}$$

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix}$$

Therefore:

$$\tilde{z}^{[l]} = W^{[l]} X + b^{[l]}$$

$$(n^{c_l}, m) = (n^{c_l}, n^{c_{l-1}})(n^{c_{l-1}}, m) + (n^{c_l}, 1)$$

This is still
a vector

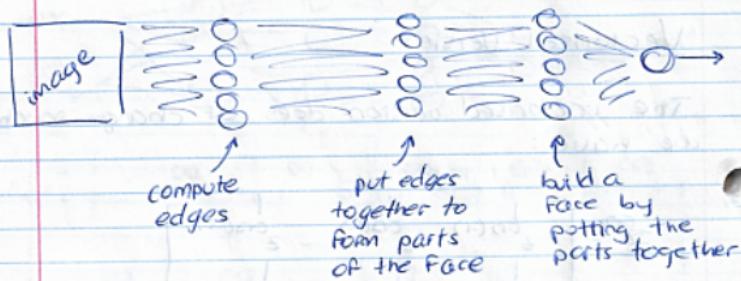
broadcasted

$$(n^{c_l}, m)$$

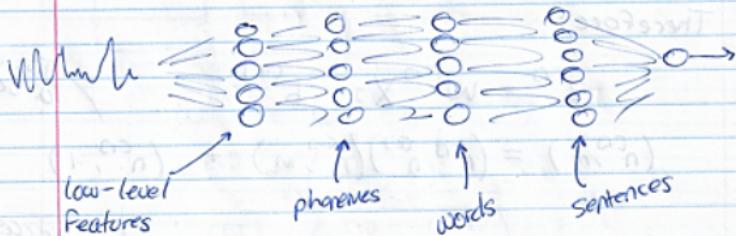
Why Deep Representations?

The idea is that further layers in the neural network can compute more and more complex functions by combining the outputs of the previous layer. Conversely, the first layers in the network compute simple functions.

Example : image recognition



Example : speech recognition

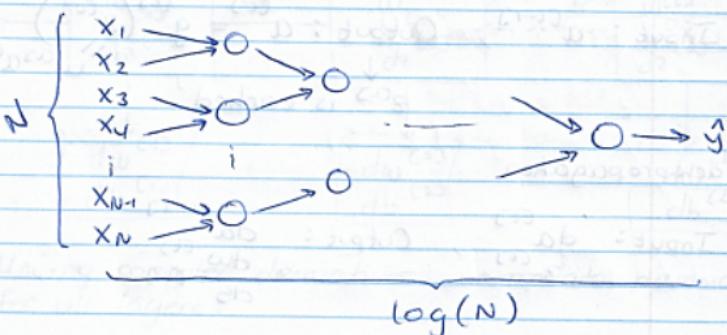


Circuit theory and deep learning

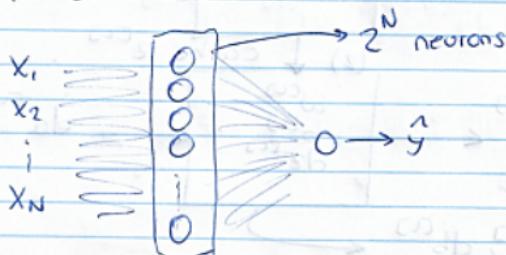
"There are functions you can compute with a 'small' L -layer deep neural network that shallower networks require exponentially more hidden units to compute".

Example: compute $x_1 \text{ XOR } x_2 \text{ XOR } \dots \text{ XOR } x_N$

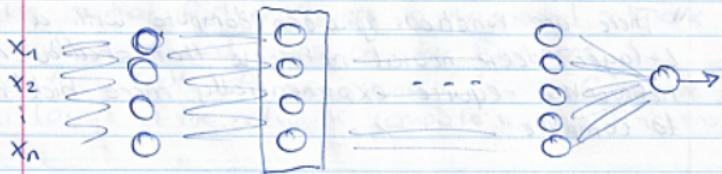
This can be done by pairing the inputs and $\log(N)$ layers:



If we attempted to do this with a single hidden layer, we would need an exponentially big number of hidden 2^N units, which would basically exhaustively enumerate all input combinations:



Building Blocks of Deep Neural Networks



Layer l : $W^{[l]}, b^{[l]}$

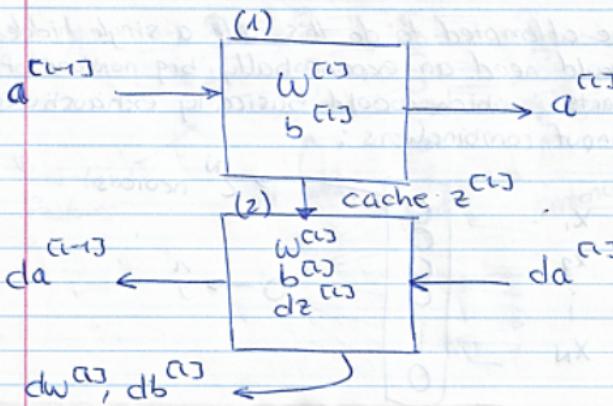
(1) Forward propagation

$$\text{Input: } a^{[l-1]}, \quad \text{Output: } a^{[l]} = g^{[l]}(z^{[l]})$$

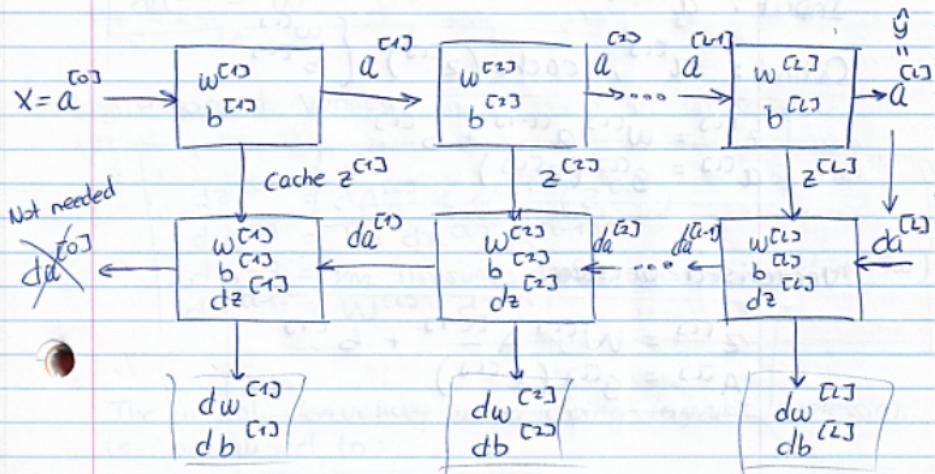
\downarrow
 $z^{[l]}$ is cached

(2) Backpropagation

$$\text{Input: } da^{[l]}, \quad \text{Output: } \begin{matrix} da^{[l-1]} \\ dw^{[l]} \\ db^{[l]} \end{matrix}$$



That was for a single layer L . To compute forward and backpropagation for the whole network, we just repeat the process for all layers.



Having computed derivatives, we update parameters for all layers:

For $L = 1 \dots L$:

$$w^{[l]} := w^{[l]} - \alpha dw^{[l]}$$

$$b^{[l]} := b^{[l]} - \alpha db^{[l]}$$

Forward and Backward Propagation

Forward propagation for layer l

Input: $a^{(l-1)}$

Output: $a^{(l)}$, cache ($z^{(l)}$) $\left\{ \begin{array}{l} w^{(l)} \\ b^{(l)} \end{array} \right.$

$$z^{(l)} = W^{(l)} a^{(l-1)} + b^{(l)}$$

$$a^{(l)} = g^{(l)}(z^{(l)})$$

Vectorised version:

$$\boxed{z^{(l)} = W^{(l)} A^{(l-1)} + b^{(l)}}$$

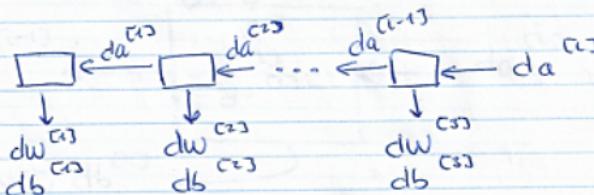
$$\boxed{A^{(l)} = g^{(l)}(z^{(l)})}$$

$$X = A^{(0)} \rightarrow \square \rightarrow \square \rightarrow \dots \rightarrow \square \rightarrow G$$

Backward propagation for layer l

Input: $da^{(l)}$

Output: $da^{(l-1)}$, $dW^{(l)}$, $db^{(l)}$



$$\begin{aligned} dz^{(i)} &= dA^{(i)} * g^{(i)}(z^{(i)}) \\ dw^{(i)} &= dz^{(i)} a^{(i-1)^T} \\ db^{(i)} &= dz^{(i)} \\ da^{(i-1)} &= W^{(i)^T} dz^{(i)} \end{aligned}$$

element-wise

Vectorised Version:

$$\begin{aligned} dz^{(i)} &= dA^{(i)} * g^{(i)}(z^{(i)}) \\ dw^{(i)} &= \text{mean}(dz^{(i)} A^{(i-1)^T}) \\ db^{(i)} &= \text{mean}(\text{np.sum}(dz^{(i)}, \text{axis}=1, \text{keepdims=True})) \\ dA^{(i-1)} &= W^{(i)^T} dz^{(i)} \end{aligned}$$

The initial derivative, when doing logistic regression, is initialised to:

$$dA^{(i)} = -\frac{y}{a} + \frac{(1-y)}{(1-a)}$$

The vectorised version $dA^{(i)}$ contains this same value but for all training examples:

$$dA^{(i)} = [dA^{(i)(0)} \quad dA^{(i)(1)} \quad \dots \quad dA^{(i)(n)}]$$

Parameters vs Hyperparameters

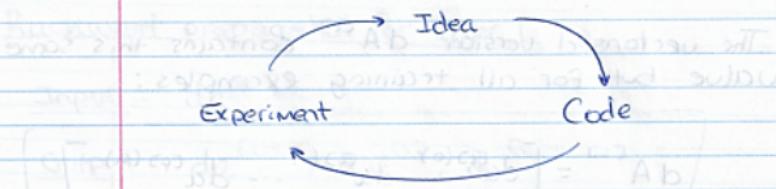
The parameters of our model are weights and biases:

$$w^{c_1}, b^{c_1}, w^{c_2}, b^{c_2}, \dots, w^{c_l}, b^{c_l}$$

Then there are other parameters that affect the learning algorithm, and that ultimately control the above parameters w, b . This last set of parameters we call hyperparameters, which include:

- Learning rate α
- # iterations
- # hidden layers L
- # hidden units $n^{c_1}, n^{c_2}, \dots, n^{c_L}$
- Choice of activation functions

Determining the hyperparameters is an empirical process. You just need to play around and see what works.



There are more systematic approaches towards determining the hyperparameters that are seen in later courses.