


why is the gradient the direction of steepest ascent?

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

Primer on partial derivatives

The question that directional derivative asks is:

what does this ridge in the direction of the vector do to the function itself?

$$\vec{v} = \begin{bmatrix} -1 \\ 2 \end{bmatrix}$$

so what we are asking is, given an arbitrary vector $\vec{a} = [a_1, a_2, \dots, a_n]^T$, we want to find what is

$$\lim_{h \rightarrow 0} \frac{f(\vec{a} + h\vec{v}) - f(\vec{a})}{h \|\vec{v}\|} := \nabla_{\vec{v}} f(\vec{a}) \text{ (definition)}$$

The definition of $\frac{\partial f}{\partial v_n}$ is

$$\lim_{h \rightarrow 0} \frac{f(\vec{a} + h\hat{i}_n) - f(\vec{a})}{h} := \frac{\partial f(\vec{a})}{\partial x_n}$$

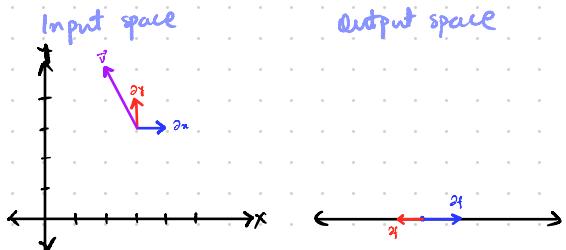
so if we have a vector $\vec{v} = [p, q]^T$, $\nabla_{\vec{v}} f(\vec{a})$ would be equal to (rate of change in n direction) $\times p +$ (rate of change in y direction) $\times q$

↳ This can be intuitively understood by looking at the graph

Therefore, $\nabla_{\vec{v}} f(\vec{a})$

$$= p \cdot \frac{\partial f(\vec{a})}{\partial x_n} + q \cdot \frac{\partial f(\vec{a})}{\partial y} = \vec{v} \cdot \vec{\nabla} f(\vec{a})$$

we define $\vec{\nabla} = [2/\partial x, 2/\partial y]^T$ = gradient operator



Partial derivative: How does change in x (or y only) changes the output?

Formally, it gives the rate of change of function as you move in the direction of that vector

We want to answer the following question: In which direction does the function f increase the fastest?

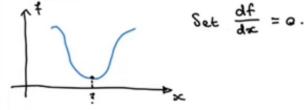
In short, we want to maximize $\nabla_{\vec{v}} f(\vec{a}) = \vec{v} \cdot \vec{\nabla} f(\vec{a})$ for some \vec{v} .

clearly, $\vec{v} \cdot \vec{\nabla} f(\vec{a})$ is maximum when \vec{v} is in the same direction as $\vec{\nabla} f(\vec{a})$.
 $\therefore \vec{\nabla} f(\vec{a})$ is the direction of steepest ascent.

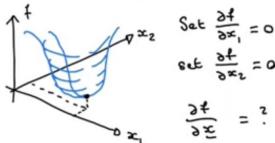
vector and Matrix derivatives

Instead of tediously doing partial derivative w.r.t. each variable, why not group all the variables into a vector and define a vector derivative, that when acted upon gives a vector when it is either at a local maxima or minima?

Main idea:
How do we find minimum of a scalar function?



And for a function of two variables?



- What if we have a function with N variables?
- Functions with intermediate variables?
- Functions producing a vector as output instead of a scalar?

Main idea:
Define vector and matrix derivatives to allow us to differentiate directly in vector/matrix form.

Definitions

- Derivative of a scalar function $f : \mathbb{R}^N \rightarrow \mathbb{R}$ with respect to vector $\mathbf{x} \in \mathbb{R}^N$:

$$\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \triangleq \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \frac{\partial f(\mathbf{x})}{\partial x_2} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_N} \end{bmatrix}$$

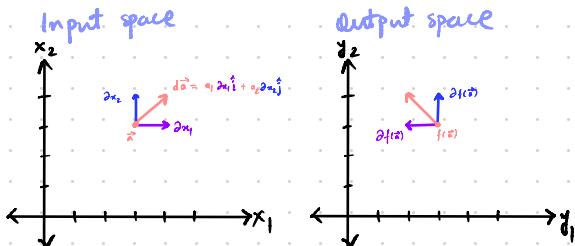
- Derivative of a vector function $\mathbf{f} : \mathbb{R}^N \rightarrow \mathbb{R}^M$ with respect to vector $\mathbf{x} \in \mathbb{R}^N$:

$$\frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} \triangleq \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \frac{\partial f_1(\mathbf{x})}{\partial x_2} & \dots & \frac{\partial f_1(\mathbf{x})}{\partial x_N} \\ \frac{\partial f_2(\mathbf{x})}{\partial x_1} & \frac{\partial f_2(\mathbf{x})}{\partial x_2} & \dots & \frac{\partial f_2(\mathbf{x})}{\partial x_N} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_M(\mathbf{x})}{\partial x_1} & \frac{\partial f_M(\mathbf{x})}{\partial x_2} & \dots & \frac{\partial f_M(\mathbf{x})}{\partial x_N} \end{bmatrix}$$

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}^2$$

$$f(x_1, x_2) = (y_1, y_2)$$

$$\frac{\partial f(x_1, x_2)}{\partial \mathbf{x}} =$$



Jensen's Inequality

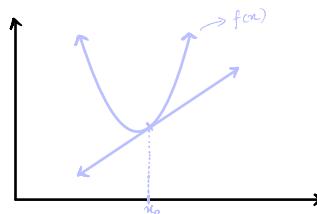
Definition: A random variable is called integrable if expected value exists and is well-defined. i.e.

$$E(X) < \infty$$

$$\Rightarrow \int x f(x) dx < \infty \text{ if } X \text{ is continuous, } f_X(x) = \text{pdf}$$

$$\Rightarrow \sum_{x \in X} x p_X(x) < \infty \text{ for discrete RV, } p_X(x) = \text{pmf.}$$

Definition: A function f is convex iff all the points lie above the tangent drawn at arbitrary point x_0 .



Jensen's Inequality: For a convex function $f: \mathbb{R} \rightarrow \mathbb{R}$, and an integrable random variable X ,

$$E(f(X)) \geq f(E(X))$$

Proof: Let the slope of the tangent at any arbitrary point of the function x_0 , be m .

The equation for the line will be : $y = m(x - x_0) + f(x_0)$

clearly, $f(x) \geq m(x - x_0) + f(x_0) \quad \forall x_0, x \in \mathbb{R}$

We replace $x \rightarrow X$ & $x_0 \rightarrow E(X)$ and get

$$f(X) \geq m(X - E(X)) + f(E(X))$$

Taking Expectation on both side and using linearity of Expectation, we get

$$\begin{aligned} E(f(X)) &\geq m(E(X) - E(X)) + E(f(E(X))) \\ &= 0 + f(E(X)) \end{aligned}$$

$$E(f(X)) \geq f(E(X))$$

Q.E.D.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

$$\begin{array}{c} \xrightarrow{x} \xrightarrow{y} \\ \downarrow \quad \downarrow \\ |0\rangle \rightarrow \frac{|0\rangle + |1\rangle}{\sqrt{2}} \end{array} \quad \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} = U$$

$$|1\rangle \rightarrow \frac{|0\rangle - |1\rangle}{\sqrt{2}} \quad U^{-1} = U^+$$

$$U^* = U^T = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}$$

$$\frac{\pi}{4} + \cos(\frac{\pi}{4}) + \cos(\frac{\pi}{4} + \cos(\frac{\pi}{4})) + \cos(\frac{\pi}{4} + \cos(\frac{\pi}{4}) + \cos(\frac{\pi}{4} + \cos(\frac{\pi}{4})))$$

$$= \frac{\pi}{2}$$

$$\det f(n) = \sin(n)$$

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} \rightarrow x_1 y_1 - (x_1 y_2 - x_2 y_1) + 2 x_2 y_2$$

$$\begin{pmatrix} y_1 \\ y_2 \end{pmatrix}, \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \rightarrow y_1 x_1 - (y_1 x_2 - y_2 x_1) + 2 y_2 x_2$$

$$\begin{pmatrix} 1 \\ 0 \end{pmatrix} \begin{pmatrix} 0 \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$\begin{pmatrix} 0 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$A = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

$$n + f'(cn) + f'(cn + f'(cn)) + f'(cn + f'(cn) + f'(cn + f'(cn)))$$

$$\begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$$

$$A^2 = I$$

$$\frac{f(w_1 x_1 + w_2 x_2)}{w_1 + w_2} \leq \frac{w_1 f(x_1) + w_2 f(x_2)}{w_1 + w_2} \quad \left((\psi_1 \vec{b}_1 + \psi_2 \vec{b}_2 + \dots + \psi_n \vec{b}_n) \right.$$

$$\left. , (\lambda_1 \vec{b}_1 + \lambda_2 \vec{b}_2 + \dots + \lambda_n \vec{b}_n) \right)$$

$$b_{ij} = \langle b_i, b_j \rangle$$

$$\sum_{i=1}^n \psi_i \alpha \vec{b}_i, \lambda_1 \vec{b}_1 + \dots + \lambda_n \vec{b}_n \rangle$$

$$\sum_{i=1}^n \psi_i \alpha \vec{b}_i, \sum_{j=1}^n \lambda_j \vec{b}_j \rangle$$

$$\sum_{i=1}^n \sum_{j=1}^n \psi_i \lambda_j \alpha \vec{b}_i, \vec{b}_j \rangle$$

$$[\psi_1, \psi_2, \dots, \psi_n] \begin{bmatrix} b_{11} & b_{12} & \dots & b_{1n} \\ \vdots & & & \vdots \\ b_{n1} & \dots & b_{nn} \end{bmatrix} \begin{bmatrix} \lambda_1 \\ \vdots \\ \lambda_n \end{bmatrix}$$

$$[\psi_1, \dots, \psi_n] \begin{bmatrix} \sum_{j=1}^n b_{ij} \lambda_j \\ \vdots \\ \sum_{j=1}^n b_{nj} \lambda_j \end{bmatrix} = \psi_1 \cdot \sum_{j=1}^n b_{1j} \lambda_j + \dots + \psi_n \sum_{j=1}^n b_{nj} \lambda_j = \sum_{i=1}^n \sum_{j=1}^n \overbrace{\psi_i b_{ij} \lambda_j}^{\psi_i \times b_{ij} \lambda_j}$$

$$\begin{aligned}A^2 &= I \\A^{-1} &= II\end{aligned}$$

七

$$e^{tn} = I + At + \frac{A^2 t^2}{2!} + \frac{A^3 t^3}{3!} + \dots$$



$$\begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$$

$$\begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix} = A_1^2 - A_1^n = 0 \quad n \geq 2$$

$$\begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} n \\ 1 \end{pmatrix} \rightarrow \begin{pmatrix} 1 \\ 0 \end{pmatrix} \rightarrow 0$$

$$= I + At + \frac{t^2}{2!} I + \frac{At^3}{3!} + \dots$$

(all terms are zero)
 $| \vec{a} | \sin \theta$

$$= I(1 + \frac{t^2}{2!} + \dots) t + At(t + \frac{At^2}{3!} + \dots)$$

$\Rightarrow \vec{r} = I T \cos \theta$

$$1 - \frac{t^2}{2!} + \frac{t^4}{4!} - \frac{t^6}{6!},$$

$$\cos(t) \quad I \cos(it) + At$$

$$\vec{r} \cdot \vec{s} = r s \cos \theta$$

$$\sin \theta = \sqrt{1 - \left(\frac{\vec{v} \cdot \vec{r}}{|\vec{v}|}\right)^2}$$

$$A_{1t} = I + A_{1t+0}$$

$$\frac{d\vec{r}}{dt} = (\vec{I} + \vec{A}(t))$$

$$I = I_0 \cos(\Omega t) + A t \sin(\Omega t)$$

$$e^{A_1 t} e^{A_2 t} = (I + A_1 t)(I + A_2 t)$$

$$= I + A_2 t + A_1 t^2$$

operator

$\nabla: \mathbb{R} \rightarrow \mathbb{R}^d$

$$f(x, y) = \begin{pmatrix} x \\ y \end{pmatrix}$$

$\checkmark f(x, y) \rightarrow \text{Vector}$

direction of
steepest

	F	S	sun	Monday	Tuesday	Wednesday	Thursday	Friday
15 th	16 th	17 ^m	18 th	19 th	20 th	21 st	22 nd	
f	K	N	N	N	N	S	S	N
	Y	S	N	N	N	S	S	N
	S	N	N	S	S	N	N	-

d_2^2 is big

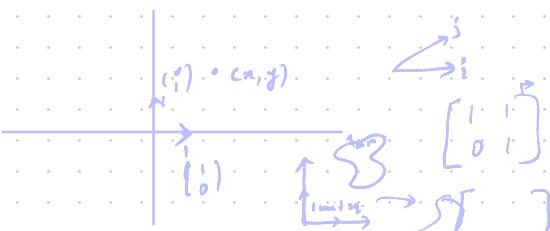
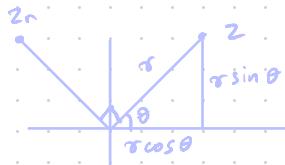
d_n

The diagram illustrates a three-phase system with three vertical lines representing phases. The top phase has an upward-pointing arrow labeled a_1 . The middle phase has a downward-pointing arrow labeled a_2 . The bottom phase has a downward-pointing arrow labeled a_3 . The resultant vector of the three phases, $a_1 + a_2 + a_3$, is shown as a horizontal line pointing to the right. This resultant vector is enclosed in a bracket labeled C_{mf} . A curved arrow at the bottom indicates the direction of rotation of the vectors.

\rightarrow select $a_1 \rightarrow$
 $\exists x y : a_2 < y < a_2$
 \rightarrow select a_2

$$\begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$$

matrices as transformational function



$$\text{rotation} \quad \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad \boxed{z_n = \bar{z} e^{\pi i / 2}}$$

$$\begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ -1 \end{bmatrix} \quad \bar{z} = r e^{i\alpha}$$

$$\begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$\frac{1}{2} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 4 \\ 3 \end{bmatrix} = \begin{bmatrix} -3 \\ 4 \end{bmatrix}$$

$$P = \frac{1}{2} |10\rangle\langle 01| + \frac{1}{2} |11\rangle\langle 11|$$

$$|\psi\rangle = \frac{1}{\sqrt{2}} |10\rangle + \frac{1}{\sqrt{2}} |11\rangle$$

$$\begin{array}{c} 4 + 3i \\ -3 + 4i \end{array}$$

$$P = \frac{1}{2} |10\rangle\langle 01| + \frac{1}{2} |11\rangle\langle 11| + |10\rangle\langle 01| + |11\rangle\langle 11|$$

$$\begin{aligned} & \frac{1}{2} + \frac{1}{2} |10\rangle\langle 01| + \frac{1}{2} |11\rangle\langle 11| \\ & \frac{1}{2} |11\rangle\langle 11| \end{aligned}$$

$$|\psi\rangle = \frac{1}{\sqrt{2}} |10\rangle + \frac{1}{\sqrt{2}} |11\rangle$$

Born
I + |10\rangle\langle 01| + |11\rangle\langle 11|

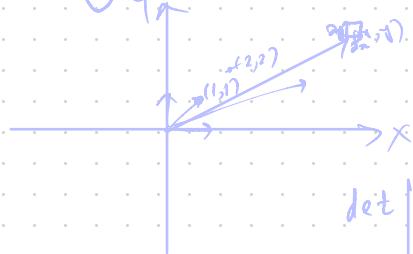
$$\frac{1}{2} (|10\rangle\langle 01| + |11\rangle\langle 11)$$

$$\begin{aligned} & \frac{1}{2} (|10\rangle\langle 01| + |10\rangle\langle 11| \\ & + |11\rangle\langle 01| + |11\rangle\langle 11|) \\ & \text{superpos. } + |10\rangle\langle 01| + |11\rangle\langle 11| \end{aligned}$$

(det)

what is a jacobian

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



$$\rho =$$

$\det(M)$ = how much
does the
matrix squish
in super.

classical / noise/
ensemble / static

$$1|10\rangle\langle 0| + \frac{1}{2}|11\rangle\langle 11|$$

det | Jacobian

$$|\psi_{AB}\rangle = \frac{1|10\rangle + 1|11\rangle}{\sqrt{2}}$$

A	B
1 0>	
1 1>	1 0> 1 1>

$$|10\rangle \rightarrow$$

$$|\psi_A\rangle = |10\rangle$$

$$|\psi_B\rangle = |11\rangle$$

$\int \int |\psi\rangle$ boundary

$$\frac{\partial f}{\partial x} \quad \frac{\partial f}{\partial y}$$

$$|10\rangle\langle 0| = \rho \rightarrow \underline{\text{pure state}}$$

$$U = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$$

n

$$\begin{bmatrix} 1 \\ 0 \end{bmatrix} |10\rangle$$

n →

$$\begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} ? \end{bmatrix}$$

$$U^{-1} = U^+ = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$$

↗

[] → ①

$$|10\rangle \rightarrow (|1+\rangle + |1-\rangle)/\sqrt{2}$$

$$|1+\rangle \rightarrow \sum_{i=1}^n \sum_{j=1}^n$$

$$\begin{vmatrix} 1-\lambda & 0 \\ 0 & -\lambda \end{vmatrix} = 0$$

$$\rho = \frac{1}{2} |10\rangle\langle 0| + \frac{1}{2} |11\rangle\langle 11|$$

$$|1, 0\rangle \rightarrow \text{degenerate} \quad \rho = |10\rangle\langle 0|$$

A = ^{spectral decomp}

$$0 \quad p = 1/2|0\rangle\langle 0| + 1/2|i\rangle\langle i|$$

≥ 2

Born

$$U_{AB} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & 1 & 0 \\ 0 & 1 & -1 & 0 \\ -1 & 0 & 0 & -1 \end{bmatrix}$$

$$K_{ij}^o = \sqrt{\lambda_i} \mathbb{I}_S \otimes |j\rangle\langle j| U_{SE} |\mathbb{I}_S \otimes |i\rangle$$

$$K_{ij}^o = \sqrt{\lambda_i}$$

Back propagation

We have a cost function $C(X)$, where $X = \text{training data}$ and X is $n \times d$ dimensional matrix containing n data points of dimension d .

In case where the dataset is MNIST, $d = 28 \times 28 = 784$ and $n = 60,000$

We define x^i as the i^{th} datapoint in the dataset. $\therefore X = [x^1 \ x^2 \ x^3 \dots x^n]^T$

and define $c(x^i)$ the cost per example/datapoint

which is $c(x^i) = \sum_{k=1}^m (y_k^i - \hat{y}_k^i)^2$, here \hat{y}_k^i is the prediction of the neural network of the datapoint x^i and y_k^i is the actual label that of the datapoint x^i .
(a vector of dimension m)
(also a vector of dimension m)

In case of MNIST, $m = 10$ since there are 10 digits we want to predict.

Now, we define $C(X) = \frac{1}{n} \sum_{i=1}^n c(x^i) = \text{average cost over all training dataset} = \text{Total cost of the network.}$

$$\therefore C(X) = \text{Total cost of the network} = \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^m (y_k^i - \hat{y}_k^i)^2$$

Here $\hat{y}_k^i = f_\theta(x^i)$, meaning f_θ is our neural network parameterized by θ that takes in input x^i gives \hat{y}_k^i as output.

Hence $\theta = \text{weights and biases.}$

In case of 3B1B MNIST example, we have

- i) 784 neurons as an input layer
- ii) 2 hidden layers each consisting of 16 neurons
- iii) Output layer consisting of 10 neurons

$$\therefore \text{Total number of weights} = \underbrace{16 \times 784}_{\substack{\text{(going from input to 1st hidden layer)} \\ \text{for 1st hidden layer neurons}}} + \underbrace{16 \times 16}_{\substack{\text{(from 1st hidden layer to 2nd hidden layer)} \\ \text{for output layer neurons}}} + \underbrace{16 \times 10}_{\substack{\text{(2nd hidden layer to output)}}} = 12,960$$

$$\text{and total bias terms} = \underbrace{16 + 16 + 10}_{\substack{\text{for 1st hidden layer neurons} \\ \text{for output layer neurons}}} = 42$$

\therefore In total, there are $12,960 + 42 = 13,002$ weights and biases

Θ is a collection of all these weights and biases

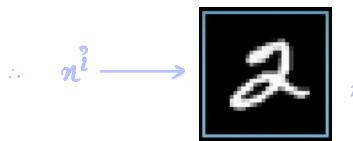
$$\text{So we can write } \hat{y}^i = f(x^i; w_0, w_1, \dots, w_{13,001})$$

$$\text{and the cost of the network as } C(X; w_0, w_1, w_2, \dots, w_{13,001})$$

we want $-\nabla_w C(X; w_0, w_1, \dots, w_{13,001})$ so that we can update the weights and biases and decrease the cost of the network.

Backpropagation is the algorithm for computing this crazy complicated gradient efficiently.

We will now be focusing our attention how the cost of the network is affected by a single arbitrary datapoint, say $x^i \rightarrow$ an image of 2 from MNIST



For an untrained network,

\hat{y}_0^i	0.5	0	0
	0.8	1	1
	0.2	2	2
	1.0	3	3
	0.4	4	4
	0.6	5	5
	1.0	6	6
	0.0	7	7
	0.2	8	8
	0.1	9	9

we wish to change the activations of the output layer such that it correctly classifies x^i as a 2.

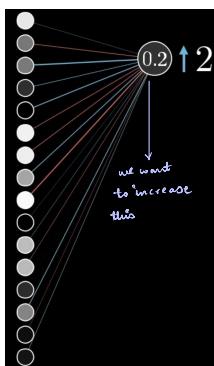
But we cannot directly change the activations, since we only have control over weights & biases but we will see it will be helpful to keep track of what changes we need to make

0.5 \downarrow 0	0
0.8 \downarrow 1	1
0.2 \uparrow 2	2
1.0 \downarrow 3	3
0.4 \downarrow 4	4
0.6 \downarrow 5	5
1.0 \downarrow 6	6
0.0 \downarrow 7	7
0.2 \downarrow 8	8
0.1 \downarrow 9	9

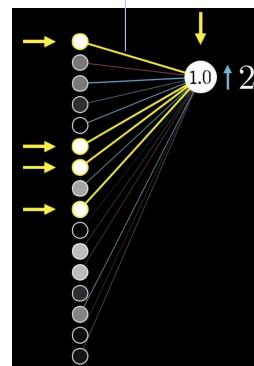
Increasing this more important than this

focusing on the neuron y_3^i as shown in the figure above* and below, we see

$$(0.2) = \sigma(w_0 a_0 + w_1 a_1 + \dots + w_{n-1} a_{n-1} + b)$$

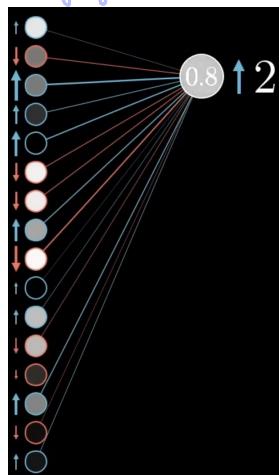


We can do that by increasing the weights in proportion to the activations of the previous layer i.e. $\Delta w_i \propto a_i$



or by changing the activations of the neurons connected to it

Clearly, we can also change the biases to get the desired value of the output neurons

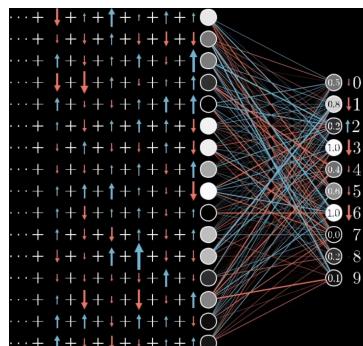


in proportion to the weights i.e. neuron with high positive weights should be increased more while those with negative weights should be decreased.

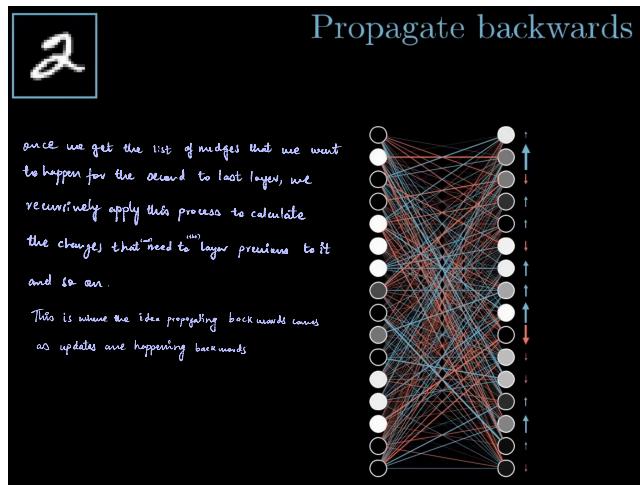
$$i.e. \Delta a_i \propto w_i$$

We obviously cannot change the activations directly, but (we will see) it helps keep track of this

we remember that the changes that we wanted to happen to the $w+B$ and activations as described above is only based on y_3^i , the 3rd neuron of our output layer, and therefore we add the relative changes based on other output neurons as well i.e. y_k^i , $1 \leq k \leq 10, k \neq 3$, in proportion to the corresponding weights and also in proportion to how much each of these neurons need to change.



After adding nudges from all the output neurons, the final nudges that we get for the second to last layer, represents the desire of y^i , and how the network should change so that it minimizes the cost function.



We remind ourselves that all of this is how a single training example y^i wishes nudge all the weights and biases. We need to consider how other training wish of nudging all the weights and biases.

							...	Average over all training data
w_0	-0.08	+0.02	-0.02	+0.11	-0.05	-0.14	...	→ -0.08
w_1	-0.11	+0.11	+0.07	+0.02	+0.09	+0.05	...	→ +0.12
w_2	-0.07	-0.04	-0.01	+0.02	+0.13	-0.15	...	→ -0.06
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
$w_{13,001}$	+0.13	+0.08	-0.06	-0.09	-0.02	+0.04	...	→ +0.04

Averaging all the training examples' wishes, we get a final list of averaged nudges over all training examples/dataset is (proportional to) -ve gradient of the cost function $C(X; \theta)$

$$-\eta \nabla C(w_1, w_2, \dots, w_{13,001}) = \begin{bmatrix} -0.08 \\ +0.12 \\ -0.06 \\ \vdots \\ +0.04 \end{bmatrix}$$

we use this gradient descent step in our gradient descent algorithm.

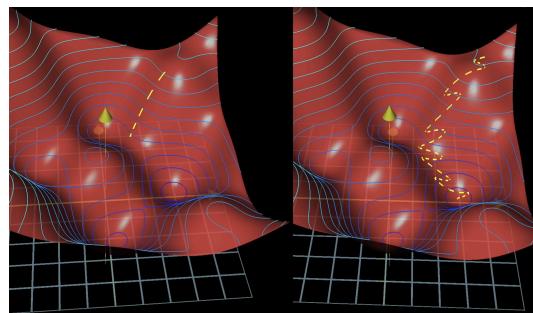
In practice, it takes computers an extremely long time to add up the influence of every single training example, every single gradient step. Therefore we will use an approximation called **stochastic gradient descent** which would speed up the process of our gradient descent algorithm.

(SGD) (This is what actually)

Stochastic Gradient Descent: We randomly shuffle our training data and divide it into "Minibatches", for example, each one having 100 training example, then we compute the gradient descent step according to the mini-batches using backpropagation. This would not be the true gradient descent step i.e. this wouldn't be the most efficient step downhill, but it each of the mini-batch does give a pretty good approximation and more importantly, it will gives us a significant computational speed-up.

If we are to plot the trajectories of our network under the relevant cost surface

Rather than a
Carefully calculating
man who takes
slow, deliberate
steps downhill



SGD would be like
a drunk man stumbling
aimlessly downhill
but taking quick steps

Summary

Backpropagation is the algorithm for determining how **a single training example** would like to nudge the weights and biases.

A true gradient descent step would involve doing this for all tens of thousands of training examples and averaging the desired changes that you get. But since doing this computationally slow, we do **SGD** to speed up the process and converge to a **local minimum**.

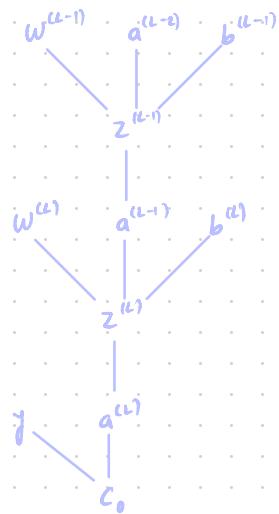
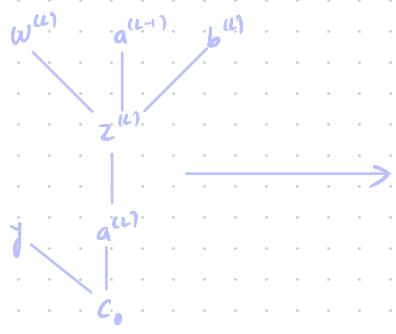
Backpropagation Calculus

$$C(\dots) = (a^{(L)} - y)^2$$

$$a^{(L)} = \sigma(w^{(L)}a^{(L-1)} + b^{(L)})$$

$$\text{def } z^{(L)} = w^{(L)}a^{(L-1)} + b^{(L)}$$

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



$$A = \sum_{i=1}^n \lambda_i |i> \langle i|$$

$$f(A) = \sum_{i=1}^n f(\lambda_i) |i> \langle i|$$

$$Y = B + X * K ; \text{ Want } \frac{\partial E}{\partial K}, \text{ Given } \frac{\partial E}{\partial Y}$$

$$\frac{\partial E}{\partial K} = \frac{\partial Y}{\partial K} \frac{\partial E}{\partial Y} = \frac{\partial}{\partial K} (B + X * K) \frac{\partial E}{\partial Y}$$

$$= (0 + X * 1) \frac{\partial E}{\partial Y} = X * \frac{\partial E}{\partial Y}$$

$$\begin{aligned}\frac{\partial E}{\partial K_{12}} &= \frac{\partial Y_1}{\partial K_{12}} \frac{\partial E}{\partial Y_1} + \frac{\partial Y_2}{\partial K_{12}} \frac{\partial E}{\partial Y_2} + \dots + \frac{\partial Y_d}{\partial K_{12}} \frac{\partial E}{\partial Y_d} \\ &= X_2 * \frac{\partial E}{\partial Y_1} + 0\end{aligned}$$

$$\therefore \frac{\partial E}{\partial K_{12}} = X_2 * \frac{\partial E}{\partial Y_1}$$

$$\boxed{\frac{\partial E}{\partial K_{ij}} = X_j * \frac{\partial E}{\partial Y_i}}$$

$$\text{we have, } Y_i = B_i + \sum_{j=1}^n X_j * K_{ij}$$

Given $\frac{\partial E}{\partial Y}$, what is $\frac{\partial E}{\partial B}$?

$$\begin{aligned}\frac{\partial E}{\partial B_2} &= \frac{\partial Y_1}{\partial B_2} \frac{\partial E}{\partial Y_1} + \frac{\partial Y_2}{\partial B_2} \frac{\partial E}{\partial Y_2} + \dots \\ &= \frac{\partial E}{\partial Y_2}\end{aligned}$$

$$\therefore \frac{\partial E}{\partial B_2} = \frac{\partial E}{\partial Y_2} \Rightarrow \frac{\partial E}{\partial B_2} = \frac{\partial E}{\partial Y_2}$$

$$\Rightarrow \boxed{\frac{\partial E}{\partial B} = \frac{\partial E}{\partial Y}}$$

Wrong way of solving the question

$$\begin{aligned}\frac{\partial E}{\partial X_3} &= \frac{\partial Y_1}{\partial X_3} \frac{\partial E}{\partial Y_1} + \frac{\partial Y_2}{\partial X_3} \frac{\partial E}{\partial Y_2} + \dots + \frac{\partial Y_d}{\partial X_3} \frac{\partial E}{\partial Y_d} \\ &= K_{13} * \frac{\partial E}{\partial Y_1} + K_{23} * \frac{\partial E}{\partial Y_2} + \dots + K_{d3} * \frac{\partial E}{\partial Y_d} \\ \frac{\partial E}{\partial X_1} &= K_{11} * \frac{\partial E}{\partial Y_1} + K_{21} * \frac{\partial E}{\partial Y_2} + \dots + K_{d1} * \frac{\partial E}{\partial Y_d}\end{aligned}$$

$$\text{Now, Given } \frac{\partial E}{\partial Y}, \text{ calculate } \frac{\partial E}{\partial X}$$

$$\text{Given } Y = B + X * K, \frac{\partial E}{\partial X} = \frac{\partial E}{\partial Y} * K$$

$$\therefore \text{Given } Y_i = B_i + \sum_{j=1}^n X_j * K_{ij}, \quad \frac{\partial E}{\partial X_j} = \frac{\partial E}{\partial Y_1} * K_{1j} + \frac{\partial E}{\partial Y_2} * K_{2j} + \dots + \frac{\partial E}{\partial Y_d} * K_{dj}$$

$$\frac{\partial E}{\partial X_j} = \sum_{i=1}^d \frac{\partial E}{\partial Y_i} * K_{ij}$$

$$\mathcal{E} = -\sum_{i=1}^n \log(y_i^{d_i^*} \cdot (1-y_i)^{1-d_i^*})$$

$$= -\log \left(\prod_{i=1}^n y_i^{d_i^*} \cdot (1-y_i)^{1-d_i^*} \right)$$

$$d_i^* = 1 \quad y_i = 0.9$$

$$0.9 \approx 1$$