

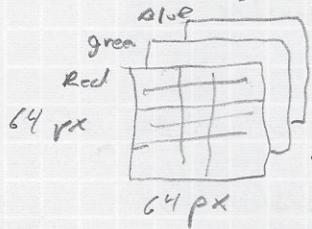
Neural Networks and Deep Learning

Week 2 Neural Networks Basics

Binary Classification

Have an input image and an output label 1 or 0, cat or noncat

e.g. Computer stores 3 matrices corresponding to the Red, Green, and Blue channels of each image.



→ unroll these into a single feature vector X

assuming $64 \text{px} \times 64 \text{px}$ image our feature vector has $64 \times 64 \times 3 = 12288$

$$n = n_x = 12288$$

$$X \rightarrow Y$$

Notation

a single pair (X, Y) $X \in \mathbb{R}^{n_x}, Y \in \{0, 1\}$

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

$$M = m_{\text{train}} = \# \text{ examples} \quad M_{\text{test}} = \# \text{ test examples}$$

to put all of the training examples into a more compact notation, we're going to take the training set inputs and stack them in columns to define matrix X

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

number of rows is equal to the number of values for each example of your input data
has m columns

This convention makes implementation in python much easier.

$$X \in \mathbb{R}^{n_x \times m} \quad X.\text{shape} = (n_x, m)$$

It is also helpful to stack your outputs

$$Y \in \mathbb{R}^{1 \times m} \quad Y.\text{shape} = (1, m)$$

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

Logistic Regression

A learning algorithm you use when the output label, Y , is either 0 or 1, so a binary classification problem.

Given input feature vector X , want \hat{y} = your estimate of Y

You want \hat{y} to be the probability of the chance that Y is equal to one given the input features X .

$$\hat{y} = P(Y=1|X)$$

$$X \in \mathbb{R}^{n_X}$$

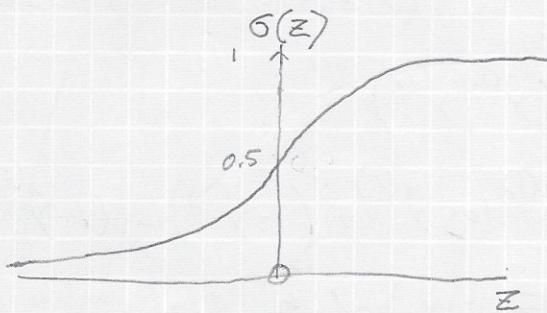
$$\text{Parameters: } w \in \mathbb{R}^{n_X}, b \in \mathbb{R}$$

(X and w are n_X dimensional vectors, b is a real number)

Output: $\hat{y} = ?$ how do you set \hat{y}

One option would be $= w^T X + b$, which is linear regression, but this isn't a good option for binary classification because you want \hat{y} ("Y-hat") to be the chance that Y is equal to one. So \hat{y} should be between 0 and 1 and it's difficult to enforce that because w transpose X can be much bigger than one or it can be negative, which doesn't make sense for probability.

Output $\hat{y} = \sigma(w^T X + b)$ ← sigmoid function



Sigmoid of z crosses the vertical axis at 0.5

Sigmoid function:

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

notice: - if z is very large e^{-z} will be close to zero,

$$\text{so } \frac{1}{1+0} = 1$$

- if z is very small e^{-z} is a large negative number so $\frac{1}{1+\text{big number}} \approx 0$

Next, in order to train the parameters w and b you need a cost function.

Logistic Regression Cost Function

Recap:

For a single training example $y^{(i)} = \sigma(w^T x^{(i)} + b)$, where $\sigma(z^{(i)}) = \frac{1}{1+e^{-z^{(i)}}}$

Given $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$, want $\hat{y}^{(i)} \approx y^{(i)}$ (want your prediction to be close to ground truth)
Superscript (i) means the i -th training example.

We could define the loss in terms of the square error, so:

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

But in logistic regression people don't usually do this because the optimization function becomes non-convex [and] where you have multiple local optimum so gradient descent might not find the global optimum.

Loss function is a measure we define to measure how good \hat{y} is when the true label is y .

In LR, this is the loss function we use which does give a convex shape that will perform well in gradient descent:

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

You want this to be as small as possible.

$$\text{If } y=1: \mathcal{L}(\hat{y}, y) = -\log \hat{y} \leftarrow (\text{because second term is zero } (1-1))$$

$$\text{want } \log \hat{y} \text{ to be large, want } \hat{y} \text{ large}$$
$$\text{If } y=0: \mathcal{L}(\hat{y}, y) = -\log(1-\hat{y}) \leftarrow \text{first term will be zero}$$

want $\log(1-\hat{y})$ to be large, want \hat{y} small

Because we are using sigmoid \hat{y} will never be greater than 1. This is one of many functions that if $\hat{y}=1$ we try to make it very large, and if zero we try to make it really small.

Loss = with respect to one training example

Cost = measures how you are doing on the entire training set (average losses)

Cost function:

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

$$\text{expanding with loss function} \quad = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)}) \right]$$

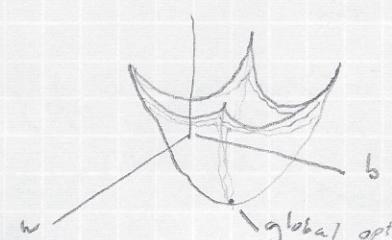
We're going to try to find parameters w and b that minimize the cost function $J(w, b)$

Gradient Descent Algorithm

Training ("learns") the parameters w, b on your training set. By finding w, b that minimize $J(w, b)$

3D plot
visualizes gradient descent with w as a single R base a single 2D

the cost function is some surface above the axes w, b



the height of the surface is the cost $J(w, b)$ at a certain point

global optimum, the minimum of cost function

so to find a good value for w, b we'll initialize them to some value. For logistic regression almost any initialization method works, usually you initialize to zero, random also works.

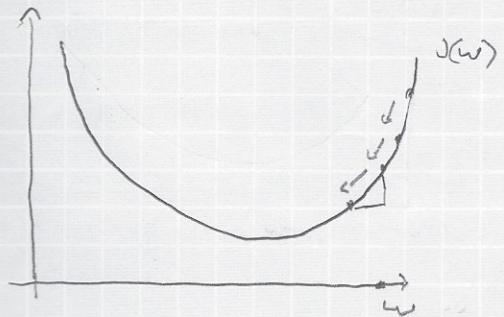
G.D. works by starting at the initialization and stepping in the direction of greatest descent, the steepest downhill direction. Then stepping further iterations down, until you converge on something close to the global optimum.

Repeat {

$$w := w - \alpha \frac{dJ(w)}{dw}$$

learning rate

in code this derivative term will be written just "dw"



just looking at w gradient decent takes a step down as per the slope (the derivative) at that point in steps per the learning rate

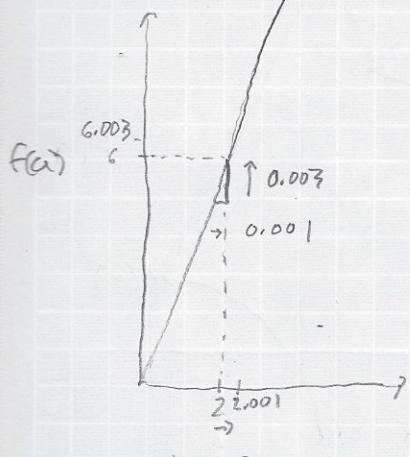
$J(w, b)$

$$w := \alpha \frac{dJ(w, b)}{dw}$$

$$b := \alpha \frac{dJ(w, b)}{db} \quad \left. \right\} "db" \text{ in code}$$

Derivatives

We want to gain an intuitive understanding of derivatives. We have a function $f(a) = 3a$, a straight line



$a = 2$	$f(a) = 6$
$a = 2.001$	$f(a) = 6.003$
$a = 5$	$f(a) = 15$
$a = 5.001$	$f(a) = 15.003$

If you nudge a over a little bit goes up three times as much so we'll say that the slope (derivative) of $f(a)$ at $a=2$ is 3, at $a=5$ also 3

slope is friendly term, derivative is scary sounding term
 height / width of our triangle $\frac{0.003}{0.001}$

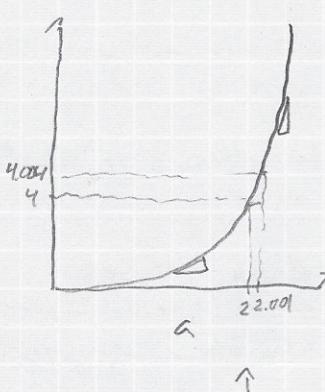
$$\frac{d f(a)}{d(a)} = 3 \quad \leftarrow \text{This says that the slope of } f(a) \text{ if you nudge } a \text{ just a little bit is 3}$$

can also write this as $\frac{d}{da} f(a)$

All this says is that if I nudge a a little bit I expect $f(a)$ to go up by three times as much

- the nudge you give in the formal definition is smaller than this 0.001, its actually infinitesimally small amount, and infinitesimally small amount.
- On a straight line, we can see that the derivative is the same everywhere, it doesn't change.

More Derivatives



$$f(a) = a^2$$

$$\begin{aligned} a &= 2 \\ a &= 2.001 \end{aligned}$$

$$\begin{aligned} f(a) &= 4 \\ f(a) &\approx 4.004 \end{aligned}$$

goes up by $4x$

$$\begin{aligned} a &= 5 \\ a &= 5.001 \end{aligned}$$

$$\begin{aligned} f(a) &= 25 \\ f(a) &\approx 25.010 \end{aligned}$$

goes up by $10x$

$$\begin{aligned} \frac{d}{da} f(a) &= 4 \quad \text{at } 2 \\ &= 10 \quad \text{at } 5 \end{aligned}$$

the ratio of the height width triangles are different at different points on the curve

Calculus textbooks list the formulas for various functions. Taking a Calc class will take you through the proofs that these are the correct derivative formulas for the respective functions. Here we're just going to give them.

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

More examples

$$f(a) = a^3 \quad \frac{d}{da} f(a) = 3a^2 \quad \begin{aligned} a &= 2 \\ a &= 2.001 \end{aligned} \quad \begin{aligned} f(a) &= 8 \\ f(a) &\approx 8.012 \end{aligned}$$

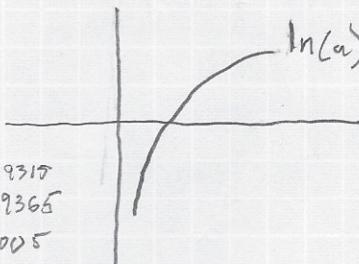
$\hookrightarrow 3 \times 2^2 = 12$

$12x$

$$\begin{aligned} f(a) &= \log_e(a) \\ &= \ln(a) \end{aligned} \quad \frac{d}{da} = \frac{1}{a}$$

$$\begin{aligned} a &= 2 \\ a &= 2.001 \\ \Delta &= 0.001 \end{aligned} \quad \begin{aligned} f(a) &= 0.69314 \\ f(a) &\approx 0.69366 \\ \Delta &= 0.0005 \end{aligned}$$

$\hookrightarrow \frac{1}{2}$ goes up by half as much



- The slope of the line varies depending on where you are for these functions.
- You can look up the formulas for the derivatives of common functions. You don't need to be able to prove them to train and use NW.

Computation Graph

Computations of a neural network are organized in terms of a forward pass (or forward propagation) in which we compute the output of the neural network, followed by a backward pass (or back propagation) step, in which we use to compute gradients (or compute derivatives).

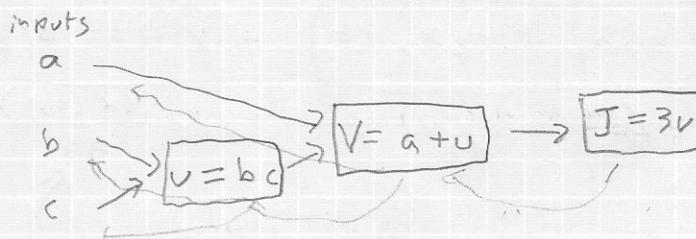
The computation graph explains why it is organized this way.

E.g. Simple example $J(a, b, c) = 3(a + bc)$

3 steps to solve

$$\begin{array}{c} J(a, b, c) = 3(a + bc) \\ \downarrow \\ u = bc \\ v = a + u \\ \downarrow \\ J = 3v \end{array}$$

Computation graph:



The computation graph comes in handy when you have some output variable (J) you want to optimize.

In logistic regression J will be the cost function.

When computing derivatives there will be a left to right pass in the direction of the arrows

Derivatives with Computation graph

$\frac{dJ}{dv} = ? = 3$ in the example above "d v "

attracts attracts

$$a \rightarrow v \rightarrow J$$

$$\frac{dJ}{da} = ? = 3 = \frac{dJ}{dv} \cdot \frac{dv}{da} \leftarrow \text{chain rule incalculus}$$

if you change a how much does it affect the final output J ?

$$\frac{dv}{da} = 1$$

by having computed $\frac{dJ}{dv}$ it helps you compute $\frac{dJ}{da}$

when coding you want a final output variable $\frac{\text{Final Output Var}}{\text{dvar}}$
you always do this so we just call it "dvar"

$$\frac{dJ}{du} = 3 = \frac{dJ}{dv} \cdot \frac{dv}{du} = d\lambda$$

$\overbrace{3}^3 \quad \overbrace{3}^3$

$$\frac{dJ}{db} = \frac{dJ}{dv} \cdot \frac{dv}{db} = 6 = d\mu$$

(chain rule in calc)

most efficient way to compute derivatives is to follow the right to left arrows and use the derivatives you have already computed to help you by

$$\frac{dJ}{dc} = \frac{dJ}{dv} \cdot \frac{dv}{dc} = 9 \quad \text{using calc's chain rule}$$

$\overbrace{3 \times 3}^3$

Logistic Regression Gradient Descent

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

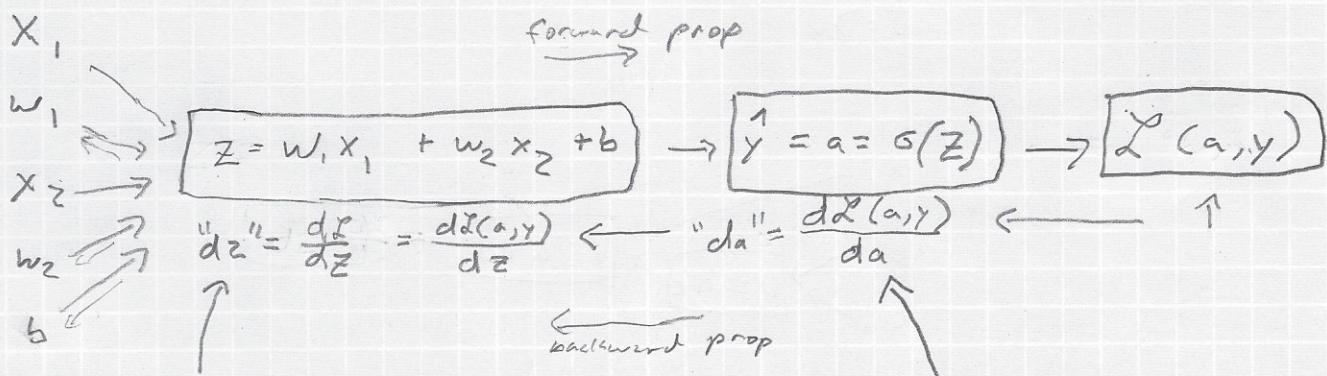
$$\hat{y} = a = \sigma(z)$$

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

*with respect to one training example

For this example let's say we have only two features X_1 and X_2

Computation graph



In L.R. we modify params w & b in order to reduce this loss.

The formulas you need to calculate the derivatives will be provided in this course. The derivative of our loss function is

$$-\frac{y}{a} + \frac{1-y}{1-a} =$$

$$= a - y$$

Again, you can derive these from the formulas above using calculus, but we're just going to look them up.

The final step in backprop is compute how much you want to change (update) w and b . Similarly, we can show that these are the equations for dw and db .

$$\frac{\partial L}{\partial w_1} = "dw_1" = x_1 \cdot dz$$

$$"dw_2" = x_2 \cdot dz$$

$$"db" = dz$$

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

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

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

Then you can perform the updates.
This is all just for one example.

Gradient Descent on m examples

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

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

$$dw_1^{(i)}, dw_2^{(i)}, db^{(i)}$$

$$(x^{(i)}, y^{(i)})$$

The above is for a single example denoted by superscript (i)

Just as the overall cost function is an average of the losses, the derivative of w_1 with respect to the cost function is also going to be the average of derivatives with respect to w_1 of the individual loss terms.

$$\overbrace{\frac{\partial}{\partial w_1} J(w, b)}^{\substack{\text{Symbol} \\ \text{just means} \\ \text{partial} \\ \text{derivative} \\ \text{because it is} \\ \text{one of several} \\ \text{terms}}} = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial \mathcal{L}(a^{(i)}, y^{(i)})}{\partial w_1}}_{dw_1^{(i)} - (x^{(i)}, y^{(i)})}$$

↑

compute these derivatives and average them

Gives you the overall gradient that you can use to implement gradient descent

To wrap this all up into an algorithm that you can use:

Initialize variables

$$J=0, dw_1=0, dw_2=0, db=0 \leftarrow \text{we will use these as accumulators}$$

use a for loop over the training set and compute derivatives with respect to each training example and then also add them up.

For $i=1$ to m :

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

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

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

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

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

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

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

$\int_{n=2}$ doing this assuming only two features
otherwise you would do this for every feature

Divide by m

$$J /= m$$

$$dw_1 /= m; dw_2 /= m; db /= m$$

So then to implement one step of gradient descent, you then

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

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

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

everything in this example would then get run again multiple times to perform multiple steps of gradient descent.

However, the algorithm above uses two For-loops which is very inefficient. Vectorization techniques help you avoid these. This is really important now in the era of big data.