

0 - 255
black white

training images

784_{px} - training image
28x28

784_{px} \Rightarrow 0, 1, 2, ..., 9
10 classes

representation as matrix

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

Each row
784 columns long

forward propagation
is taking a picture and
get probability out
but we need good weights and

biases to
make these
predictions
& the whole
point of
machine learning
is that we'll
LEARN these
weights and
biases so we
can run an algorithm
to optimize these weights
and biases and that's
called backpropagation

unactivated first layer = we get z_1 by multiplying
input layer x with weights and bias

input layer x
no processing yet

forward propagation

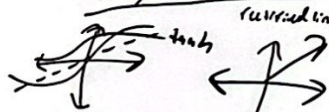
$$A^{(2)} = x (784 \times m)$$

$$z^{(3)} = w^{(2)} A^{(2)} + b^{(3)}$$

$$A^{(1)} = g(z^{(2)}) = \text{ReLU}(z^{(2)})$$

take image run through matrix
& compare what output is going to be

APPLY ACTIVATION FUNCTION



ReLU(x)

$$= \begin{cases} x & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases}$$

$A^{(1)} = \text{ReLU}(z^{(2)})$
apply to every
value of $z^{(2)}$

$$z^{(2)} = w^{(1)} A^{(0)} + b^{(2)}$$

$$A^{(2)} = \text{Softmax}(z^{(2)})$$

each of the
at the output
the softmax
activation are
... 12 +

the weight
function in
each layer
and it's a
weight is
multiplied
by another
constant
bias term

BACKWARDS
PROPAGATION

going opposite
way and starting
with prediction
and going to
find out how much
our prediction
deviated by the
actual label
basically we
error and then
where going
to see and
then where going
to see how much
each of the
previous weights
and biases contributed
to that error and
then we adjust
these things
accordingly

derivative
of last function
with respect to
weights in layer 2

$$\begin{aligned} W^{(3)} &= W^{(3)} - \alpha \delta w^{(3)} \\ b^{(3)} &= b^{(3)} - \alpha \delta b^{(3)} \\ W^{(2)} &= W^{(2)} - \alpha \delta w \\ b^{(2)} &= b^{(2)} - \alpha \delta b \end{aligned}$$

Update
Weights

$$\delta z^{(2)} = A^{(2)} - y$$

$$\delta w^{(2)} = \frac{1}{m} \sum \delta z^{(2)}$$

$$\delta b^{(2)} = \frac{1}{m} \sum \delta z^{(2)}$$

$$\delta z^{(1)} = w^{(2)T} \delta z^{(2)}$$

$$\delta w^{(1)} = \frac{1}{m} \sum \delta z^{(1)}$$

$$\delta b^{(1)} = \frac{1}{m} \sum \delta z^{(1)}$$

get predictions
and subtract
actual labels
from it

derivative of
the activation
function because
we have to undo
the activation
function to
get the proper
error for the
layer

After we do all this math
finding our weights and biases with
forward and back ward propagation
and figure out each weight term and
bias contributed to the error we
update our parameters accordingly

then we do it - and
it's a cycle of con-
tinuous improvement so that the
prediction is so close to
the actual thing

Soft max
activation function
Each of 10 nodes correspond
to each of the 10 digits
that could be recognized
we want each of them to have
a probability

Once you apply tanh
or sigmoid function to all
the nodes in a layer
there's no longer linear
once you move on to the second
layer that's now adding a
second layer of complexity &
non linear complexity rather
than just a linear model

Unactivated
second layer
nodes
each of the
at the output
the softmax
activation are
... 12 +

