

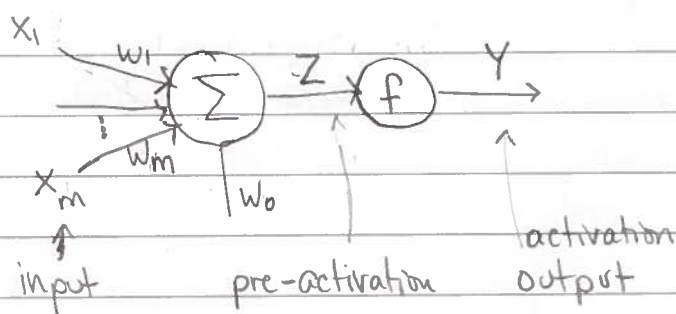
Neural networks

View 1: stochastic gradient descent for classification and regression with a potentially very rich hypothesis class

View 2: brain inspired network of neuron-like computing elements learn distributed representations

View 3: a method for building applications that make predictions based on huge amounts of data in very complex domains.

Basic element: "neuron" "unit" "node"



$$y = f\left(\sum_{j=1}^m x_j w_j + w_0\right) = f(W^T x + w_0)$$

f : "activation function"

w : "weights"

w_0 : "threshold" "bias"

differentiable!

Training

Given • a loss function $L(\text{guess}, \text{actual})$

• dataset $\{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$

Do (stochastic) gradient descent, adjusting w to minimize

$$\sum_i L(\text{NN}(x^{(i)}; w), y^{(i)})$$

(We will discuss regularization next time.)

We already know two versions

classifier w/ hinge loss

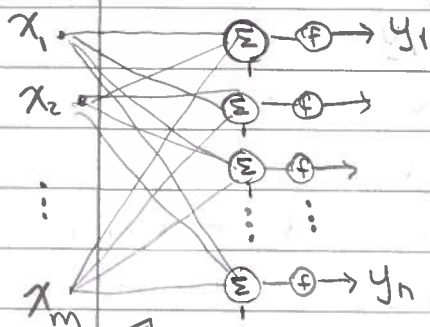
regressor w/ quadratic loss

$f = \text{identity}$

L4-2

Networks!

A layer is several units, not connected to each other.



"fully connected" if
all the inputs to each
unit are the same

n units

n outputs

m inputs

W is an $m \times n$ matrix

W_0 is an $n \times 1$ column vec

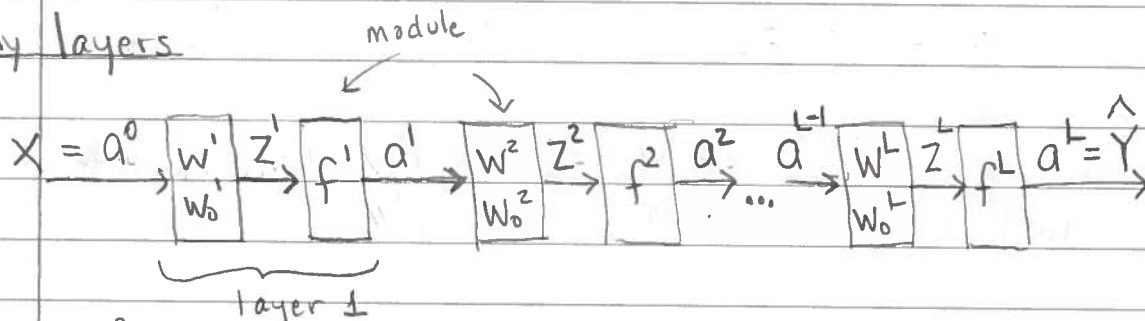
X is an $m \times 1$ column vec

Y is an $n \times 1$ column vec

$$y = f(W^T X + W_0)$$

\uparrow
 $n \times 1$ $n \times 1$
 apply elementwise

Many layers



layer l

m^l inputs

$n^l = m^{l+1}$ outputs

$W^l : m^l \times n^l$

$W_0^l : n^l \times 1$

f^l : activation function

$$Z^l = W^{lT} a^{l-1} + W_0^l \quad (n^l \times 1)$$

$$a^l = f^l(Z^l) \quad (n^l \times 1)$$

What the f!

What if we just let f be the identity? {Leaving off
no for effect}

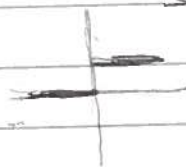
$$Y = a^L = W^{L^T} a^{L-1} = W^{L^T} W^{L-1^T} \dots W^{1^T} X$$

So

$$Y = W^{\text{total}} X$$

Having all those layers did not change the representational capacity of the network: non-linearity is crucial.

Choices for f



$$\text{step}(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{otherwise} \end{cases}$$

derivative is everywhere
0 or undefined



$$\text{ReLU}(z) = \begin{cases} 0 & \text{if } z < 0 \\ z & \text{otherwise} \end{cases}$$

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

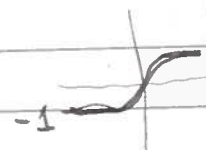
"rectified
linear unit"



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

sigmoid
logistic

interpretable
as a probability



$$\tanh(z)$$

Could be anything, but these are the common ones,
esp ReLU in internal ("hidden") layers
and σ on output.

$$\text{Softmax}(Z) = \begin{bmatrix} e^{z_1} / \sum_i e^{z_i} \\ \vdots \\ e^{z_n} / \sum_i e^{z_i} \end{bmatrix}$$

interpretable
as a
probability dist'n

↑
1-hot encoding

L4-4

Loss for classification : NLL

An alternative to hinge loss which

- is okay without regularization (but we may still need it)
- extends nicely to multi-class

Assume a softmax output layer

$$P\left(\begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix}, \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}\right) = a_i \text{ where } y_i = 1$$

prob dist activations one-hot training label product will be 1 when $y^i = 0$
 a_i when $y^i = 1$

probability assigned to the correct output

One view is that we would like to maximize the probability our network assigns to the whole dataset

$$\prod_{j=1}^{\text{dataset size}} \prod_{i=1}^{\text{output dim}} a_i^{(j) y_i^{(j)}}$$

We can maximize the log of this function instead because log is monotonic. Log is

$$\sum_{j=1}^{\text{dataset size}} \sum_{i=1}^{\text{output dim}} y_i^{(j)} \log a_i^{(j)}$$

Hints at a loss function (which we want to minimize)

$$\text{NLL}(a, y) = - \sum_{i=1}^{\text{output dim}} y_i \log a_i$$

negative log likelihood
 log loss
 cross entropy

L4-6

Modules with weights have to be able to compute
input, $\partial L / \partial \text{output} \rightarrow \partial L / \partial \text{weights}$

Training

Initializing W is important

- Don't make magnitude too big for sigmoid / tanh activation fns
- Don't set to 0
 - no gradient for relu
 - generally better to randomize for symmetry breaking
- Rule of thumb:

$$W_{ij}^l \sim \text{Gaussian}(0, 1/m^l)$$

\uparrow drawn from prob dist'n \uparrow 0 mean \uparrow variance

Magnitude of Z will be independent of # of inputs

$$W_{ij}^l \sim G(0, 1/m^l) \quad W_{0j}^l \sim G(0, 1) \quad \leftarrow \text{or } W_{0j}^l = 0$$

for $t = 1$ to T

$i = \text{random draw from } \{1 \dots n\}$ \leftarrow # of data pts

$$a^0 = x^{(i)}$$

for l from 1 to L :

$$z^l = W^{lT} a^{l-1} + W_0^l$$

$$a^l = f^l(z^l)$$

$$\text{loss} = L(a^L, y^{(i)})$$

$$dL/dz^L = (\partial L / \partial a^L) (\partial a^L / \partial z^L)$$

$$dL/dW^L = dL/dz^L \cdot dz/dW^L$$

for l from $L-1$ down to 1:

$$dL/dA^l = dL/dz^{l+1} \cdot dz^{l+1}/dA^l$$

$$dL/dZ^l = dL/dA^l \cdot dA^l/dZ^l$$

$$dL/dW^l = dL/dZ^l \cdot dz/dW^l$$

$$dL/dW_0^l = dL/dZ^l \cdot dz/dW_0^l$$

$$W^l = W^l - \eta(t) dL/dW^l$$

$$W_0^l = W_0^l - \eta(t) dL/dW_0^l$$

backprop

weight grad

sgd

* More about this later

$$\begin{cases} dL/dW^L = dL/dz^L \cdot dz/dW^L \\ W^L = W^L - \eta(t) dL/dW^L \\ W_0^L = W_0^L - \eta(t) dL/dW_0^L \end{cases}$$

L4-5

Error backpropagation : backprop

To do SGD, focus on loss with respect to 1 (x, y)

Need to compute

$$\nabla_W \text{Loss}(\text{NN}(x; W), y)$$

↑ all the W^L, W_0^L

Seems terrifying, but actually easy and cool!

The chain rule is our hero

We will grossly abuse calc notation

Start with W^L

$$\frac{\partial \text{Loss}}{\partial W^L} = \frac{\partial \text{Loss}}{\partial a^L} \cdot \frac{\partial a^L}{\partial z^L} \cdot \frac{\partial z^L}{\partial W^L}$$

depends on
loss function

↑ $f^{L'}$

depends on
 a^{L-1}

$$\frac{\partial \text{Loss}}{\partial W^1} = \frac{\partial \text{Loss}}{\partial a^L} \cdot \frac{\partial a^L}{\partial z^L} \cdot \frac{\partial z^L}{\partial a^{L-1}} \cdot \frac{\partial a^{L-1}}{\partial z^{L-1}} \cdots \frac{\partial a^2}{\partial z^2} \cdot \frac{\partial z^2}{\partial a^1} \cdot \frac{\partial a^1}{\partial W^1}$$

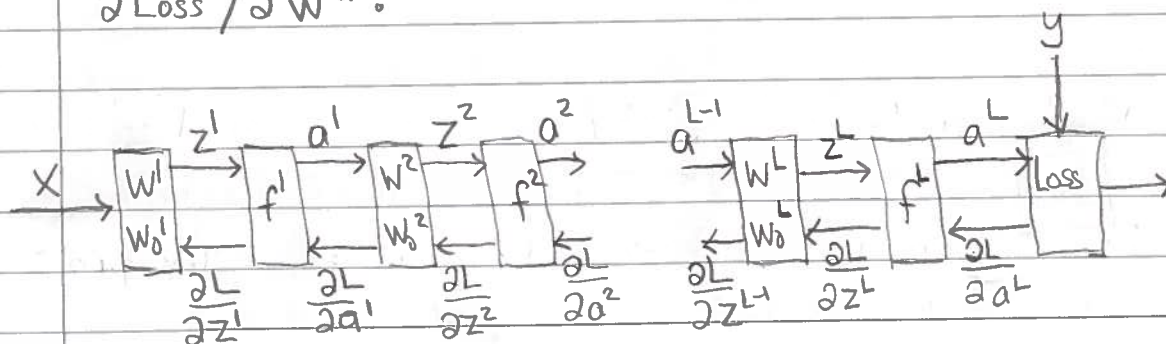
$\frac{\partial \text{Loss}}{\partial z^2}$

$\frac{\partial \text{Loss}}{\partial a^1}$

So, we can do a systematic backward propagation of errors

(I like to think of it as blame) to compute all the

$\frac{\partial \text{Loss}}{\partial W^L}$:



Every module has to be able to compute

forward : input \rightarrow output

backward : input, output, $\frac{\partial L}{\partial \text{output}} \rightarrow \frac{\partial L}{\partial \text{input}}$

L 4-7

(do NLL here)

Loss functions make assumptions about the type of input they're getting - so in our code we assume that f^L (activation function of last layer) and L (loss) go together:

Loss	f^L
squared	linear
hinge	linear
NLL ₁	sigmoid (in homework)
NLL	softmax

It turns out that directly computing

$\frac{\partial L}{\partial z^L}$ is usually easier than separately doing

$$\frac{\partial L}{\partial a^L} \text{ and } \frac{\partial a^L}{\partial z^L}$$

especially for NLL. So we will ask our software implementation of a loss function to provide a backward method that computes $\partial L / \partial z$ directly.

The shape of things

What is $\partial \text{Loss} / \partial W^L$?

A matrix of the same shape as W^L ($m^L \times n^L$)
with entries: $\partial \text{Loss} / \partial W_{ij}^L$

Similarly:

- $\partial \text{Loss} / \partial a^L$ is a vector ($n^L \times 1$) with entries $\partial \text{Loss} / \partial a_i^L$
- $\partial \text{Loss} / \partial z^L$ is a vector ($n^L \times 1$) with entries $\partial \text{Loss} / \partial z_i^L$