

* Sigmoid neuron → building block of Deep Neural Networks.

↓
or the foundational block of DNN.

↓
not limited to sigmoid, can be something like ReLU neuron or any other non-linearity

↓
generalised case: sigmoid neuron.

↓
till now i/p's: real i/p's

task: biclassification

model: linear (can't operate with non-linear data)

↓
weights and bias were tuneable (perceptron)

indicator fn.
or 1-0 loss

← loss: $\sum_i 1_{\{y_i \neq \hat{y}_i\}}$ or MSE

training: perceptron learning algo.
or brute search (MP)

evaluation: accuracy

Note: o/p plotted along with the features would result in a kind of a ramp function (for perceptron).

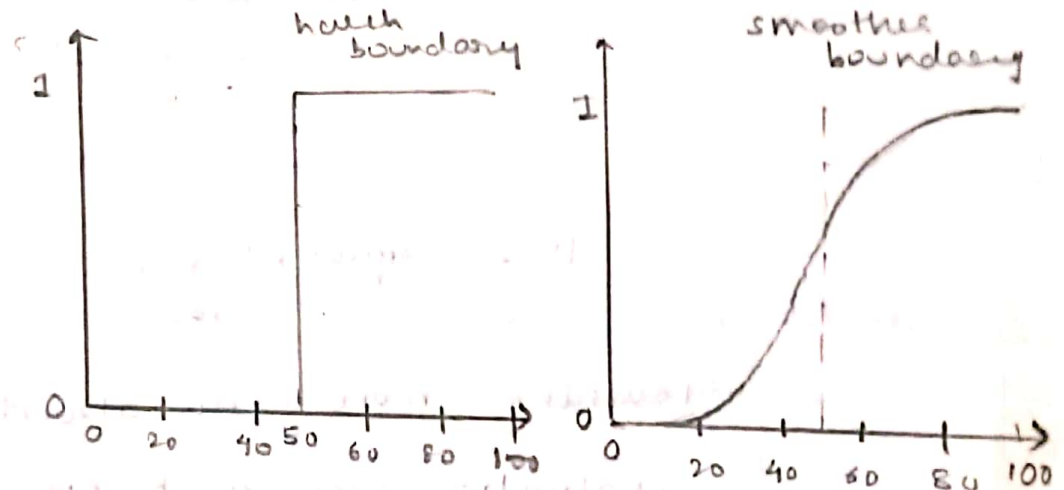
① if linearly separable)

↓
also the decision boundary is very sharp. for perceptron, rather practically the decisions are much smoother

•) a person with salary 50.1k can afford a car but 49.9k can't?

↓
boundary is harsh, and we would rather want a more smooth boundary, where there is some uncertainty associated with a range of i/ps.

salary	buy?
80	1
20	0
65	1
15	0
30	0
49	0
51	1
87	1



therefore person with a salary near to 50k also has a prob. of being able to afford a car

the probability on the right of 50 is now ~~low~~ approaching 1 and prob. on the left is not equal to 0.

↓
which is logical practically.

↓
These probabilities are continuous, and hence can also be employed with regression ~~in some way~~, and are not just constrained to classification.

↓
real valued o/p.s.

for sigmoid neuron: i/p: real

task: classification/regression

↓
real o/p ✓

model: smooth at boundaries
non-linear

(can deal with non linear
fns to a limit, only first
step in that direction)

loss: squared error
loss

↓
network of
neuron

learning: more generic algorithm

evaluation: accuracy + RMSE

(root mean square
error)

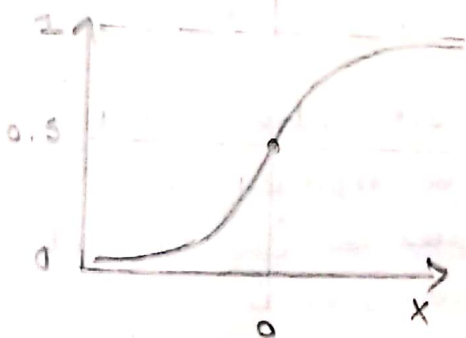
↓
since this would be
used for regression

parameters are $[w], b$

↑
the mapping b/w y and x

→ approximated using sigmoid

(*) model: a smoother fn. → sigmoid family of fns.



substitute values in fn. to
visualize.

$$y = \frac{1}{1 + e^{-(wx+b)}}$$

→ single feature

↓
logistic function

↓
always +ve

↓
between 0 + 1 (probability)

3. evaluate $\hat{\pi}(s, a)$ → single sample
 4. evaluate policy gradient
- $\nabla_{\theta} J(\theta) \approx \nabla_{\theta} \log \pi_{\theta}(s, a) \cdot \hat{\pi}(s, a)$

- *) hue / color can be used to visualize probability smoothness for two feature data, rather than plotting 3D graphs, (with certainty along the ends and uncertainty in the middle)
- *) The power of the exponential is the aggregation of the features with their weights.

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

→ scalar quantity.

- *) In perceptron model, there is no distinction between the points very close to the decision boundary and the points very far to the decision boundary as the points far from the decision boundary are very sure to be classified correctly, but the ones close would / should still have some uncertainty associated with them.

↓

can be attributed to the harsh boundary and the o/p being either 1 or 0.

↓

The sigmoid neuron deals with this issue as it assigns a ~~higher to much~~ ~~probability to~~ higher level of confidence to points far from the decision boundary as compared to the ones close to it.

- *) o/p of sigmoid can also be interpreted as probability.
- *) does not separate the i/p, rather giving a more graded o/p.

and we probably need to use parameter

* variation according to w and b for the model

$w \rightarrow -ve$, then the slope of the curve is $-ve$, as the denominator \uparrow , and total value decreases as $w \uparrow$

more $-ve$, more steeper slope of curve.

$+ve$, reverse happens

$$\frac{1}{1+e^{-w}} = \frac{1}{1+\frac{1}{e^w}}$$

$$\text{or } \frac{1}{1+e^w}$$

$$\frac{1}{1+\frac{1}{e^w}} > \frac{1}{1+e^w}$$

⑥ \rightarrow decreasing the value pushes the curve to the right

$$y = \frac{1}{1+e^{-(wx+b)}} = \left(\frac{1}{2}\right) \rightarrow \text{at the middle}$$

$$\Rightarrow e^{-(wx+b)} = 1$$

$$\Rightarrow wx+b=0$$

$$\Rightarrow x = \left(-\frac{b}{w}\right) \rightarrow \text{as } b \text{ decreases, with constant } w, x \uparrow \text{ (hence shift to the right)}$$

4. evaluate policy π

$$\pi(\theta) \approx \nabla_{\theta} \log \pi$$

$$15) \cdot \hat{\pi}(\theta) > 1$$

→ real x/p. feature

* Data and Task

(.) sigmoid neuron can be used for binary classification as well by highlighting which class the feature vector is more closer to. i.e. we are now getting a softer o/p.

↓
to take a decision, some threshold can be applied to these values. ~~to decide~~

(.) regression: given a feature vector, fitting a real value of y to that based on training data.

↓
for sigmoid, we are regressing the probability.

↓
could be used to o/p the rating for the phone (normalised btw 0 and 1)

* Loss function:

$$\sum_{i=1}^N (y - \hat{y})^2 \rightarrow \text{squared error loss fn}$$

true prediction

↓
can also be used with binary labels.

↓
for this case the point closer to the labels contribute less to the loss whereas the points away from it contribute more.

* learning algorithm

objective: give the data, we want to generalise the mapping between x and y , using sigmoid function whose parameters are w and b

- i) initialize w, b
- ii) iterate over the dataset, and compute loss on predicted value
- iii) update w and b , and repeat

need to be tuned to get a appropriate approximation

→ the generalised roadmap.

w and b need to be updated and approximated in order to minimize the loss over the dataset.

a) method 1: guesswork, : initialize $w=0, b=0$,
(by visualizing the points and the function)

straight line
($y=0.5$)

and accordingly taking the decisions to vary w and b .

change $w=1, b=0$,
somewhat better (the slope)

going from $w=1$ to $w=3$
(slope more steeper)

but, yes not scalable to higher dimensions.

moving the fn. to the right
(decreasing b)

guesswork or random changes would take a time to actually approximate

again changing w and b
a little works.

in other words, for each iteration,
 $w = w + \Delta w$
 $b = b + \Delta b$

3. evaluate $\hat{\pi}(S, a)$
4. evaluate policy gradient $\nabla_{\theta} J(\theta) \approx \nabla_{\theta} \log \pi(S, a) \cdot \hat{\pi}(S, a)$

* for guesswork were we guided by the loss function?

yes in a ~~weak~~ way the guesswork was guided by the loss function, but it can only get us so far.

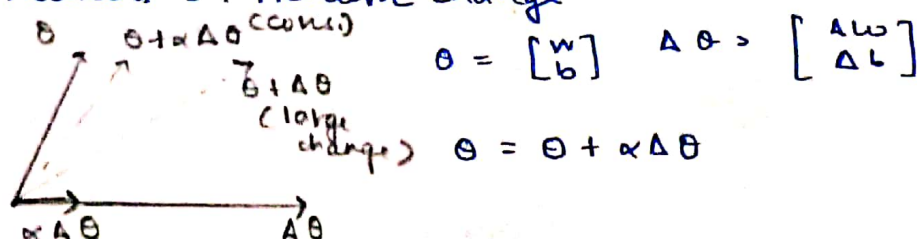
(it is easier when there are only 2 or 3 features, but for high dimensional data, this would be next to impossible to determine)

observation: Also after one point of decreasing the loss, the loss began to increase, and we had to move back a little to reach the optimal w and b for minimal loss.

for guesswork, the trajectory of optimizing w and b to minimize loss has been very random and full of fluctuations, and in high end dimensionality, where the guesswork will move in random directions and hence reaching the minima will require a large no. of steps.

so this motivates the need for a more principled algorithm that allows us to move such a way that there is gradual decrease in the loss associated with w and b .

b) method 2: iterative change



how to compute $\Delta \theta$?

↓
change or diff. vector

↓
gradient of loss with respect to θ

↓
frameworks like TensorFlow and pyTorch
have functions that give these gradients
~~and~~ and these are used to tune θ
at each iteration.

* Taylor series .

$$f(x + \Delta x) = f(x) + f'(x) \Delta x + \frac{1}{2!} f''(x) \Delta x^2 + \frac{1}{3!} f'''(x) \Delta x^3 + \dots$$

↓
value in close proximity
to x , given by
RHS

↓
for our problem, we take a small
step away from w and we need
that the loss after taking the
small step decreases (ideally)

↓
Require a relation between the loss
before and after the small step.

↓
can be accomplished using Taylor
series.

↓
where the value after the small step
is equal to the value of $f(x)$
plus some quantity, which comprises of Δx .

↓
if Δx ~~some~~ is picked such that the added
quantity is negative then it can be concluded

that the loss after the step is less than before.

↓
the more -ve, the better, since the loss will decrease by more.

↓
along with Δx and the constants, there are n^{th} order derivatives in the quantity.

↓
loss depends on b as well $\therefore b \rightarrow b + \eta b$
 $L(b) > L(b + \eta b)$

(prediction depends on both w and b)

↓
which affects the calculated loss

function of both, i.e.

$$L(w, b) > L(w + \eta \Delta w, b + \eta \Delta b)$$

$$\begin{bmatrix} w \\ b \end{bmatrix} + \eta \begin{bmatrix} \Delta w \\ \Delta b \end{bmatrix}$$

↓
moving to the ' θ ' notation

$$L(\theta) > L(\theta + \eta \theta) \quad \therefore \theta = [w, b]$$

↓ θ Taylor series for a vector

~~original eq~~

$$L(\theta + \eta u) = L(\theta) + \eta * u^T \nabla_{\theta} L(\theta) + \frac{\eta^2}{2!} * u^T \nabla^2 L(\theta) u + \frac{\eta^3}{3!} * \dots$$

$$\Delta \theta = u$$

change vector

↓
needs to be found in order to make loss less after the step.

↓
ignoring the higher order terms ($\eta \rightarrow$ very small)

$$L(\theta + \eta u) \approx L(\theta) + \eta * u^T \nabla_{\theta} L(\theta)$$

first order derivative of a vector

needs to be found in order to make this term -ve

$\nabla_{\theta} L(\theta) \rightarrow$ gradient of function which depends on θ vector wrt θ

↓

partial derivatives for all specific members.

↓

$$\begin{bmatrix} \partial f / \partial w \\ \partial f / \partial b \end{bmatrix}$$

↓

finally, we get vector of partial derivatives

↓

$L(\theta) \rightarrow$ Real number

$L(\theta + \eta u) \rightarrow$ Real number

$\eta \rightarrow$ scalar / real value

corresponding to learning rate

↓

$u^T : [\Delta w \ \Delta b]$ (vector)

$$\nabla_{\theta} L(\theta) : \begin{bmatrix} \frac{\partial L}{\partial w} \\ \frac{\partial L}{\partial b} \end{bmatrix}$$

product will also be a real no.

↓

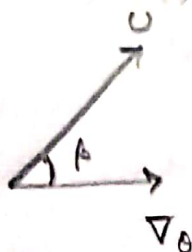
Finding u^T such that the dot product becomes -ve.

↓ Now

$$L(\theta + \eta u) - L(\theta) = \eta * u^T \nabla_{\theta} L(\theta)$$

↓ Here

$$u : \begin{bmatrix} \Delta w \\ \Delta b \end{bmatrix} \cdot \nabla_{\theta} L(\theta) : \begin{bmatrix} \frac{\partial L}{\partial w} \\ \frac{\partial L}{\partial b} \end{bmatrix}$$



$$\cos \phi = \frac{(\vec{u} \cdot \vec{\nabla_{\theta}})}{||u|| ||\nabla_{\theta}||}$$

→ this needs to be less than 0. → cos to be less than 0

→ always true

The value $\nabla^T \cdot \nabla_b L(\theta)$ will be maximum ~~when~~
~~they~~ when they are in opposite direction to each
 other, where $\cos \beta = -1$, and $\beta = 180^\circ$.



Hence v_f should be chosen in such a way that
 it is in a direction opposite to that of the
 gradient.



Gradient Descent Rule.

↓ formally

$$\begin{bmatrix} w_{t+1} \\ b_{t+1} \end{bmatrix} \rightarrow \theta_{t+1}$$

$$w_{t+1} = w_t - \eta \Delta w_t$$

$$b_{t+1} = b_t - \eta \Delta b_t$$

$$\begin{bmatrix} \Delta w_t \\ \Delta b_t \end{bmatrix} \rightarrow \nabla_{\theta} L(\theta)$$

$$\text{where } \Delta w_t = \frac{\partial L(w, b)}{\partial w} \text{ at } w = w_t \text{ and } b = b_t$$

$$\text{and } \Delta b_t = \frac{\partial L(w, b)}{\partial b} \text{ at } w = w_t \text{ and } b = b_t$$



the gradients can be computed using autograd
 in frameworks like pytorch or TensorFlow.



The algorithm involves iterating over the dataset,
 calculating the predictions, corresponding loss,
 and \therefore loss gradient and then updating the
 parameters according to the update rule.



Computing partial derivatives,

$$L = \frac{1}{n} \sum_{i=1}^n \frac{1}{2} (f(x_i) - y_i)^2$$

for convenience

$$\frac{\partial L}{\partial w} = \frac{\partial}{\partial w} \left[\frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2 \right]$$

(sumrule)

$$\Delta w = \frac{\partial L}{\partial w} = \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial w} \frac{1}{2} (f(x_i) - y_i)^2$$

(sum of individual derivatives)

•) solving for 1 out of the n terms of the sum, and then generalising.

$$\begin{aligned}
 \nabla \omega &= \frac{\partial}{\partial \omega} \left[\frac{1}{2} * (f(x) - y)^2 \right] \\
 &= \frac{1}{2} * 2 (f(x) - y) * \frac{\partial}{\partial \omega} (f(x) - y) \quad \text{chain rule} \\
 &= (f(x) - y) * \frac{\partial}{\partial \omega} (f(x)) \quad \downarrow \text{independent of } \omega \\
 &= (f(x) - y) * \frac{\partial}{\partial \omega} \left(\frac{1}{1 + e^{-(\omega x + b)}} \right) \quad \dots (i)
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial f(x)}{\partial \omega} &= \frac{\partial}{\partial \omega} \left(\frac{1}{1 + e^{-(\omega x + b)}} \right) \\
 &= \frac{-1}{(1 + e^{-(\omega x + b)})^2} \cdot e^{-(\omega x + b)} \cdot (-x) \quad \downarrow \text{chain rule} \\
 &= \frac{1}{1 + e^{-(\omega x + b)}} \cdot \frac{e^{-(\omega x + b)}}{1 + e^{-(\omega x + b)}} \cdot x \\
 &= f(x) \cdot (1 - f(x)) \cdot x
 \end{aligned}$$

replacing this value in (i), we get

$$\nabla \omega = (f(x) - y) (f(x)) (1 - f(x)) \cdot x$$

↓
computed n no.
of times, and
summed over

true o/p for the data point

↓ similarly for b ($x=1$)

$$\Delta b = (f(x) - y) (f(x)) (1 - f(x))$$

•) generalising the gradient computation over > 2 parameters



x_{ij} : notation to mark the i th data point
and j th feature in that data point



$$z = \sum_{j=1}^m w_j \cdot x_{ij}, \quad m: \text{no. of features,} \\ \text{because } \downarrow \\ m \text{ parameters.}$$



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



computing dependency of the loss on a
particular feature.



$$\Delta w_j = \sum_{i=1}^m (\hat{y} - y) * \hat{y} * (1 - \hat{y}) * x_{ij}$$

*) Evaluation : metric to check the performance of the model.



RMSE : root mean square error

(typically used for
regression problems)



$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y - \hat{y})^2}$$

and accuracy can be done in case of classification
tasks.