

5 Multilayer Perceptron

Deep Network - Simpler version

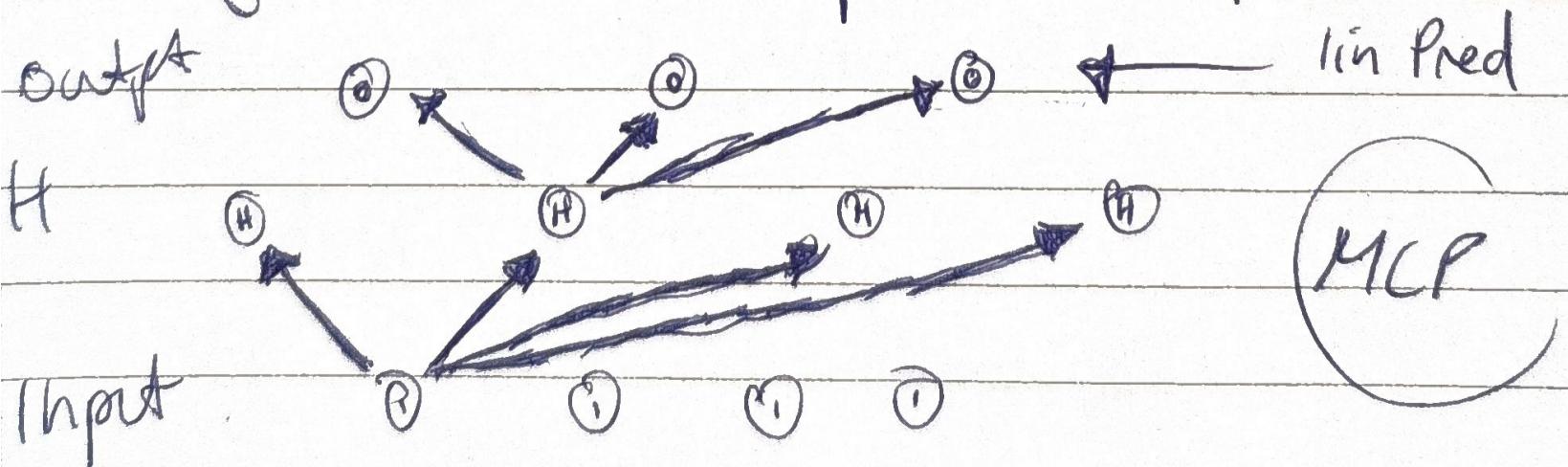
limitations of linear models is that linear assumption implies monotonicity assumption

- = increase in x must always increase output, & vice versa for decrease

NN uses a hidden layer to learn joint relationship between dimensions

Hidden layers → 

= layers between inputs & outputs



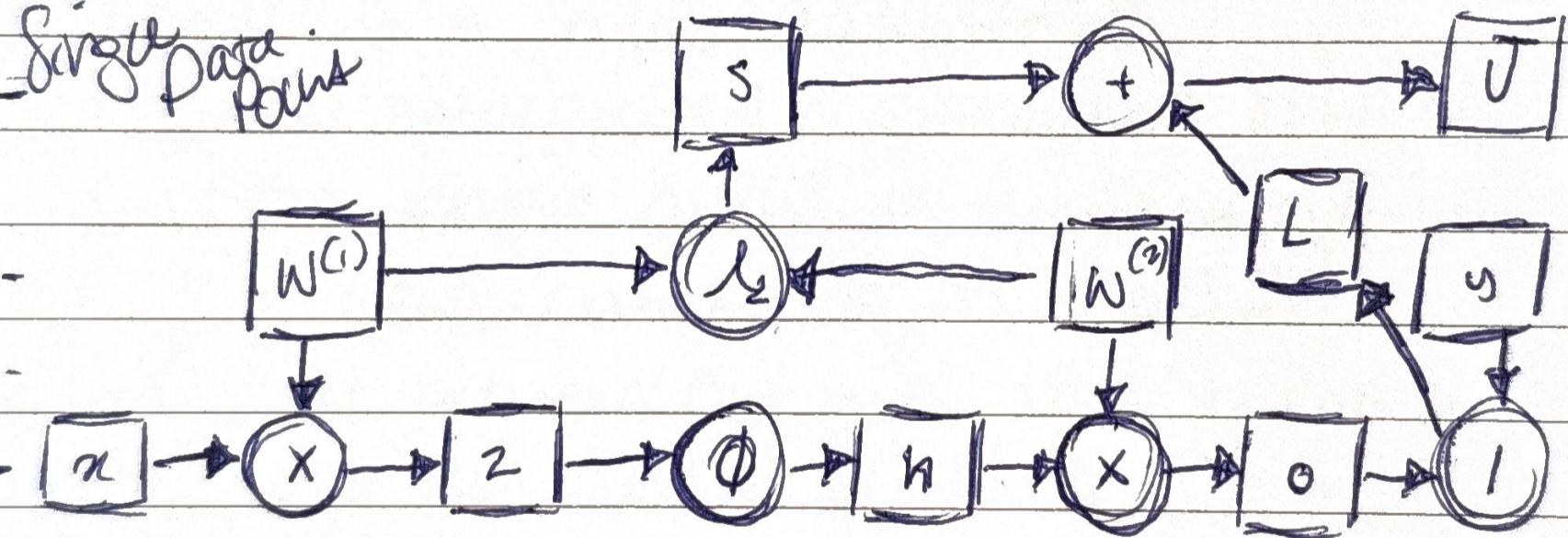
forward Pass / Propagation

(2)

Note this section is "incomplete" as there are later steps (Backpropagation) that come back and update the model

that said, forward pass calculates the output of a network given its inputs, layers, weights & transforms

single Data point



English Explanation or
other side

forward Prop - written Explanation

- We have a data point x
- We want to determine the weights $w^{(1)}$ for a linear model given $z = w^{(1)}x$ (no bias)
- These weights are obtained by minimising an objective function using gradient descent
- Using the weight/Model we calc an intermediate point $z \rightarrow z$
- These new data points are then non-linearly transformed using $\phi(z)$
- Resulting a new dataset ϕ
- A new linear model is constructed $o = w^{(2)}H$ w/ the weights calced w/ the same optim method
- An a new dataset ϕ is created
- The diff between ϕ & y is used to create the loss function value $L = L(\phi, y)$
- To prevent overfit, a reg term is calced which penalized large weights $S = \frac{\lambda}{2} \|w^{(1)}\|_F^2 + \|w^{(2)}\|_F^2$
- final objective function given by $J = L + S$
- later we will min/optimize this (J) using backpropagation

(4)

forward Prop Equations

$$\text{intermediate value } z = \mathbf{w}^{(1)} \cdot \mathbf{x}$$

$$\text{hidden active function } h = \phi(z)$$

$$\text{output layer value } o = \mathbf{w}^{(2)} \cdot h$$

$$\text{loss func } L = \lambda(o, y)$$

$$\text{reg value } S = \frac{\lambda}{2} \left(\|\mathbf{w}'\|_F^2 + \|\mathbf{w}_F\|^2 \right)$$

where λ = learning hyperparameter

$$\text{Objective func } L + S$$

(5)

Backpropagation

= the method of calc the gradient
of neural network parameters

Traverses the network in reverse
order, from output layer to
input layer, according to the
calculus chain rule

Algo store intermediate variables partial
Derivatives while calc'ng the gradient
w/ respect to some paras

- Assum funcs: $Y = f(x)$ & $Z = g(Y)$
- where input/outputs ~~are~~ (x, Y, Z) are tensor

using chain rules, compute derivative of Z
(fout) w/ respect to X (first input)

$$Z = g\left(\underbrace{f(x)}_y\right) \stackrel{\text{def}}{=} \frac{\partial Z}{\partial X} = \text{Prod}\left(\frac{\partial Z}{\partial Y}, \frac{\partial Y}{\partial X}\right)$$

(6)

Recall, chain rules

- if y depends on u which depends on x
- chain rules tells us calc Derive of y w/ respect to x
- $\frac{dy}{dx} = (\frac{dy}{du}) * (\frac{du}{dx})$
- how small change in x affects y

With vectors, we use partial derivatives to capture how each component of one vector affects another. Use MatM to comb

Chain rule for tensors:

- tensors are multi dimension arrays
Matrix = 2-D tensor
- instead of MM uses tensor contract to comb PPs
- Notation becomes + complex but concept of the chain rule remains

Prod (·) encapsulates the notation we will use to chain calcs

(7)

Recall the paras of a simple Network are given by $w^{(1)}$ & $w^{(2)}$

Obj of Back propagation is to calculate the gradients $\partial J / \partial w^{(1)}$ & $\partial J / \partial w^{(2)}$

To accomplish this, we apply chain rule and calculate

In turn, the gradient of each intermed variable & parameter

the order of calc is opposite to forward prop

Step 1 calc grads of $J = L + S$
w/ respect to loss func & reg term

$$\frac{\partial J}{\partial L} = 1 \quad \text{and} \quad \frac{\partial J}{\partial S} = 1$$

$X \rightarrow h \rightarrow O$

⑥

Step²

Compute gradient of obj func
w/ respect to variable of the
output layer O

$$\frac{\partial J}{\partial O} = \text{Prop}\left(\underbrace{\frac{\partial J}{\partial L}}_{\text{calced in}}, \frac{\partial L}{\partial O}\right)$$

calced in
 $\text{Step 1} = 1$

Step 3

calc gradient of reg term w/
respect to paras (weights)

Note; this step is done for comp^u easy &
reusability. Reg doesn't have X (data) term so
easy to use

$$\frac{\partial S_1}{\partial w^1} = \lambda w^{(1)} \quad \frac{\partial S}{\partial w^2} = \lambda w^{(2)}$$

↓ ↓

used for h

used for O

$X \rightarrow h \rightarrow O$

$$h = f(x) \quad o = g(h) \quad L = l(o, y) \quad S = \omega^1 \omega^2$$

(1) (2) (3) (4)

$$J = L + S$$

Step 4 calc gradient of object w /
respect to $\omega^{(2)}$ (or whatever
 n is closest to output o)

$$\begin{aligned} \frac{\partial J}{\partial \omega^{(2)}} &= \text{Prod}\left(\frac{\partial J}{\partial o}, \frac{\partial o}{\partial \omega^{(2)}}\right) + \text{Prod}\left(\frac{\partial J}{\partial s}, \frac{\partial s}{\partial \omega^2}\right) \\ &= \frac{\partial J}{\partial o} h^T + \lambda \omega^{(2)} \end{aligned}$$

recall ω^2 connected directly to $o = g(h)$ so
it's easier to work w/

Step 5 Obtain gradient w/ resp to $\omega^{(1)}$
need to continue down network as
 $\omega^{(1)}$ is connected to h // note we need $\frac{\partial o}{\partial h}$ first

Gradient w/ respect to hidden layer

By
 $\frac{\partial J}{\partial h}$ two parts = how much J change w/ res to o
 how much o changes w/ resp to h

$$\frac{\partial J}{\partial h} = \text{Prod}\left(\frac{\partial J}{\partial o}, \frac{\partial o}{\partial h}\right) = L^{(2)} + \frac{\partial J}{\partial o}$$

Step 6 need to get from h to z
which was the active func

Transform was element wise - hence
need to back elementwise use ⑥

$$\frac{\partial J}{\partial h} = \text{Prod}\left(\frac{\partial J}{\partial h}, \frac{\partial h}{\partial z}\right) = \frac{\partial J}{\partial h} \odot \phi'(z)$$

Step 7 now we can finally obtain
gradient $\frac{\partial J}{\partial w^{(1)}}$

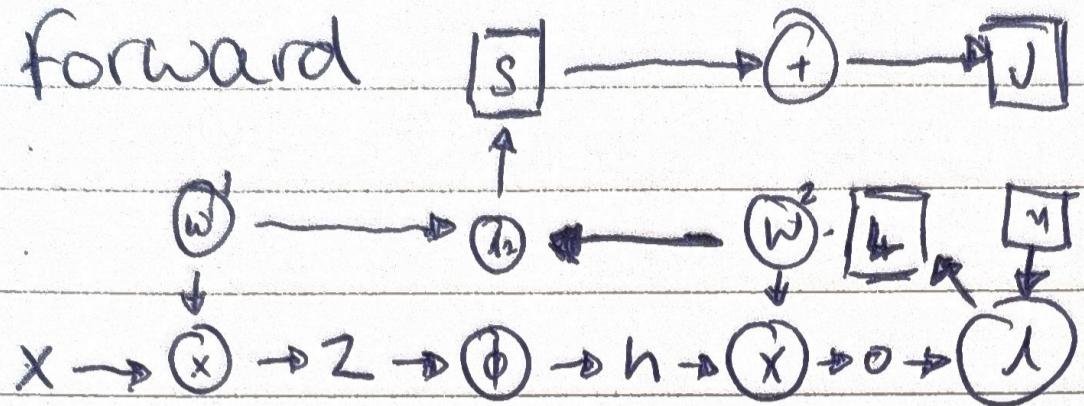
$$\begin{aligned} \frac{\partial J}{\partial w^{(1)}} &= \text{Prod}\left(\frac{\partial J}{\partial z}, \frac{\partial z}{\partial w^{(1)}}\right) + \left(\frac{\partial J}{\partial s}, \frac{\partial s}{\partial w^{(1)}}\right) \\ &= \frac{\partial J}{\partial z} X^T + \lambda w^{(1)} \end{aligned}$$

11

Graphs

X h o

forward



Backward

- ① $J = L + S \stackrel{\partial J}{\leftarrow}$ gradient of ∂J w.r.t loss ^{+ reg}
- ② $\frac{\partial J}{\partial o}$ - Diff loss func to get 0
- ③ Reg term w.r.t weights
- ④ $\frac{\partial J}{\partial w^2}$ Diff w^2 term
- ⑤ $\frac{\partial J}{\partial h}$ Diff from 0 to λ
- ⑥ undo activation for $\frac{\partial J}{\partial z}$
- ⑦ $\frac{\partial J}{\partial w^1}$ final diff w^1 term Starts from 2

(12)

Training NN

When training, forward & Backward propagation depend on each other

FP traverses graph & computes vars
 BP picks up order & optimizes weights
 Based on the opt Obj func

Note, there is a very crucial point, missed in the previous notes

forward propagation makes no optimization of the layers weights. these are set randomly

BP does all the work. FP is merely for calculating the layers output, i.e., loss, jacc & obj based on the current weights

$f = \text{calc}$

$B = \text{learn}$

(13)

NN is an iterative process. the
FP \rightarrow BP \rightarrow FP takes place many times.

each iteration is called an epoch

within the epochs the train set is often ~~split~~
split into batches & the FP/BP
steps performed for each batch

epoches decided by user (hyperpara)