

1-11-2021

week # 11 lecture # 14
Data Processing.

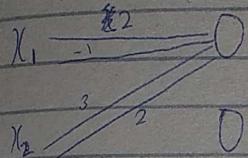
Tuesday

Activation function:

① Maxout function:

$$\max(w_1x_1+b_1, w_2x_2+b_2)$$

e.g

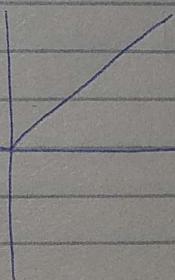


for each input there are two weights:

a) weights are doubled, memory is required more

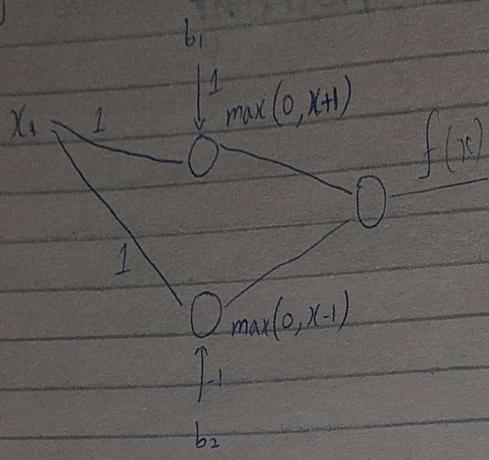
Training time is increased.

③ ReLU function:



how does this introduce non-linearity in our NN?

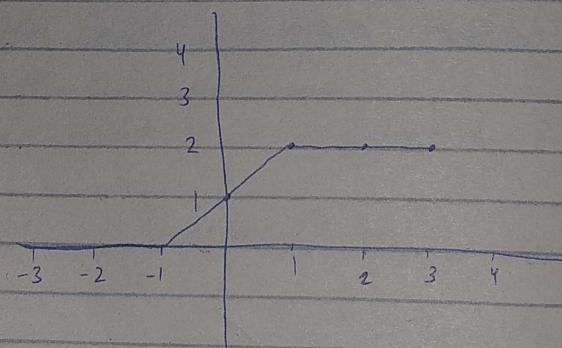
e.g



so,

$$f(x) = \max(0, x+1) - \max(0, x-1)$$

Now if we plot this graph,



$$\text{if } x=1, y=2$$

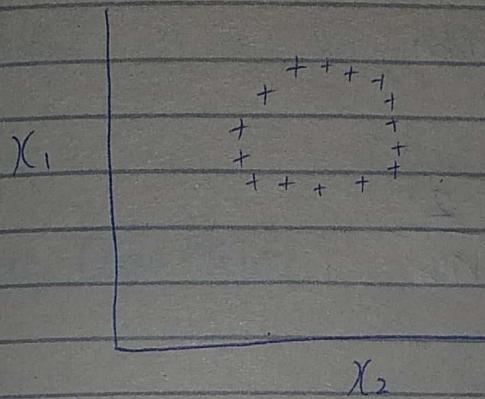
$$\text{if } x=0, y=1$$

$$\text{if } x=2, y=2$$

$f(x = -1, y = 0)$

$f(x = -2, y = 0)$

Another example:



x_1 & x_2 2 features

for this type of Dataset, a circle equation meaning
not linear function would be better for this Dataset

That's why we use ReLU function. ReLU supports
more complex function learning.

Weights Initialization:

2) Glorot Initialization.

Instead of dividing by \sqrt{n} we do.

$$w = np.random.rand(n, D) * \sqrt{\frac{2}{n_{in} + n_{out}}}$$

n_{in} : no. of input neurons

n_{out} : no. of output neurons

Researchers tell us this below one
for ReLU, we use

$$\sqrt{\frac{2}{n_{in}}}$$

or divide by

$$\sqrt{\frac{n_{in}}{2}}$$

n_{in} = no. of inputs

Above & below are same just rearranged

To overcome overfitting:-

① we did L_1 & L_2 regularization

So our loss formula was

$\frac{1}{2}$ added bcz to cancel 2 from gradient calculation

$$L = \sum_{i=1}^n L_i + \left[\frac{1}{2} \lambda \sum |w|^2 \right]$$

L_2 regularization
weights are diffused

$$L = \sum_{i=1}^n L_i + \lambda \sum |w|$$

L_1 regularization
weights are sparse

If we apply both, then it is called Elastic Net Regularization.

$$L = \sum_i l_i + \lambda_1 \sum |w| + \frac{1}{2} \lambda_2 \sum |w^2|$$

Sometimes people set different lambda's values for different ~~the~~ weights.

Use L₂ Regularization & Dataset Should be Large.

② Max Norm Constraint:-

$$|w| \leq C, \quad C = \text{constant}$$

when |w| becomes greater than C, we set that

weight value to C

usually its value is 3.0, 2.0

we are penalizing our weight values.

③ Increase Dataset size to overcome overfitting.

④ Dropout:

To ensure that our whole NN neurons are contributing

to our training, we don't have neurons that are contributing very little & others are contributing a lot.

x_1 $d w_i$ x_2

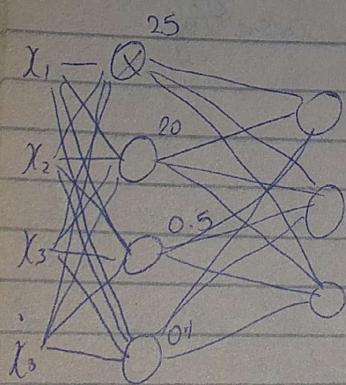
$$w \cdot w - \alpha d w_i$$

 x_3 x_4

Mini batch gradient Descent

So, concept is for different batches we disable different neurons, we randomly disable them for one batch then ~~Inactive neurons~~ would send 0 output to other neurons in forward pass.

Then what it does is other neurons which had smaller contribution before now ~~will be updated to~~ will be updated to their contribute more. & weights are updated more.



Its code:-

def forward(x):

$$\text{prob} = 0.5$$

$$h1 = \text{np. maximum}(0, W_1 \cdot \text{dot}(x))$$

$$d = \text{np. random.rand}(h1.\text{shape}) \quad \text{prob}$$

$$h2 = \text{np. maximum}(0, W_2 \cdot \text{dot}(h1))$$

$$h3 = \text{np. maximum}(0, W_3 \cdot \text{dot}(h2))$$

out =

so if forward have returned value

$$1 \ 0.3 < 0.5 \quad \checkmark$$

$$1 \ 0.4 < 0.5 \quad \checkmark$$

$$0 \ 0.7 < 0.5 \quad \times$$

$$0 \ 0.8 < 0.5 \quad \times$$

Mask > d

so first 2 neurons will be active, last 2 neurons

will be inactive & their output was multiplied by

Mask where there is 1 it will be passed forward
else, hence where there is 0 it will not be passed forward.

$$h1 = h1 * d$$

$$h2 = \text{np. maximum}(0, W_2 \cdot \text{dot}(h1))$$

At test time, no neurons would inactive, so output precessing

keep-probability was 0.5
is required because in training output was 50% bce
50% neurons was inactive,

so in test time we also want output to be

same as in train time so in test time we

multiply output of neural network ^{layers} with

keep-probability then our output would be 50%

If all neurons were active, then our output
would be doubled.

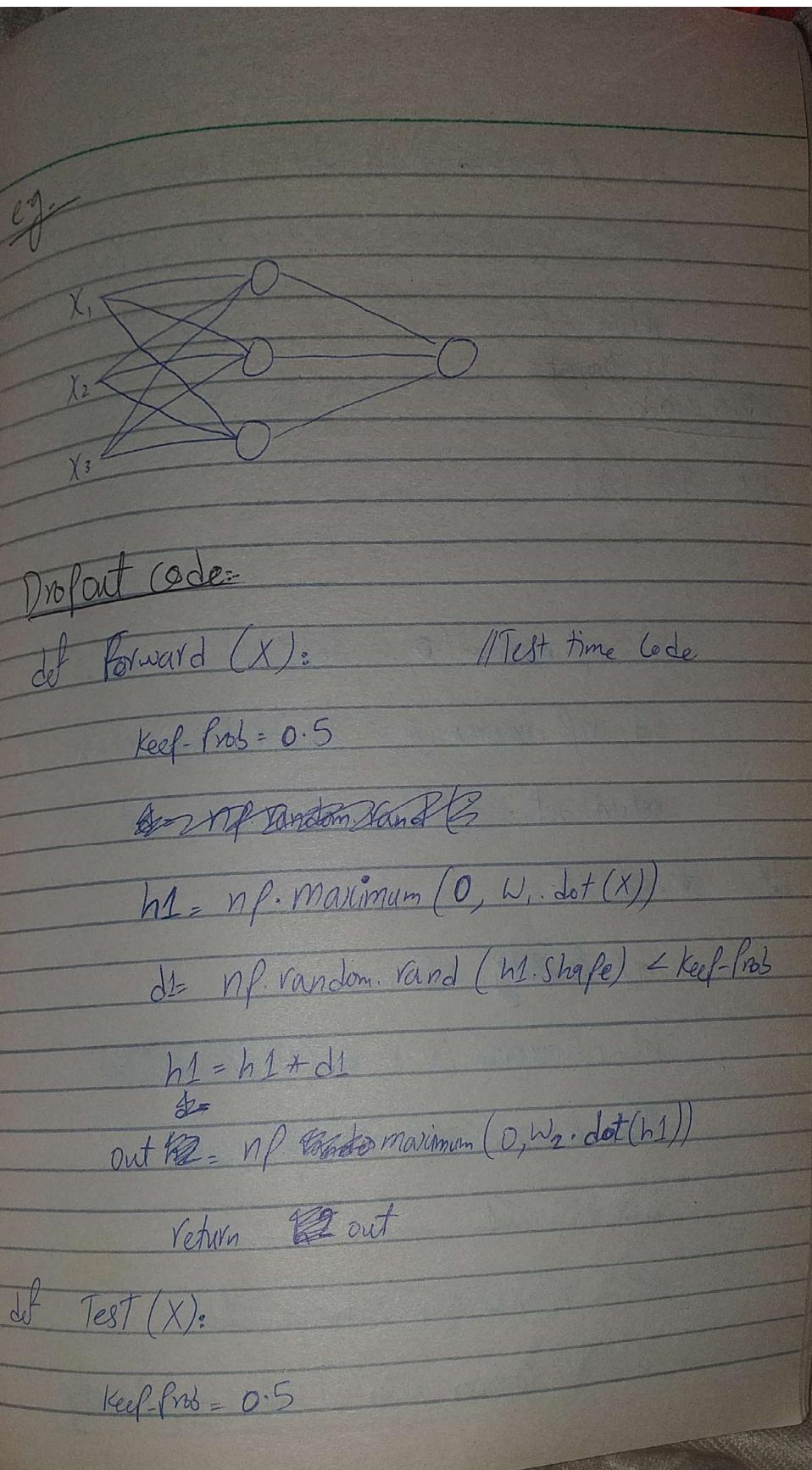
$$h_1 = np \cdot \max(0, w \cdot \text{dot}(x)) * \text{keep-probability}$$

if we
problem:-

When we are multiplying prob with every layer output, then
our test time increases. if our NN is big

Solution:- 5) Inverted Dropout.

In train time, we divide outputs of layers with
probability. So we don't have to do anything in test
time.



$h1 = np.maximum(0, w_1 \cdot \text{dot}(x)) * \text{keep_probability}$

$\text{out} = np.maximum(0, w_2 \cdot \text{dot}(h1)) * \text{keep_probability}$

return out

~~Inverted Dropout~~
~~Forward Code~~

def Test(x):

$h1 = np.maximum(0, w_1 \cdot \text{dot}(x))$

$\text{out} = np.maximum(0, w_2 \cdot \text{dot}(h1))$

return out

def forward(x):

$\text{keep_prob} = 0.5$

$h1 = np.maximum(0, w_1 \cdot \text{dot}(x))$ ~~if keep_prob~~

~~if np.random.rand(h1.shape) < keep_prob~~

~~h1 *= keep_prob~~

~~out = np.maximum(0, w_2 \cdot \text{dot}(h1))~~

$d1 = np.random.rand(h1.shape) < \text{keep_prob}$

0
0
0

$$h1 = h1 * (\alpha / \text{keep_prob})$$

$$\text{out} = \text{np.maximum}(0, w_2 \cdot \text{dot}(h1))$$

return out

Libraries for ML/DL:-

① Sklearn

② Pytorch

③ keras

④ Tensorflow

before next class check these libraries.