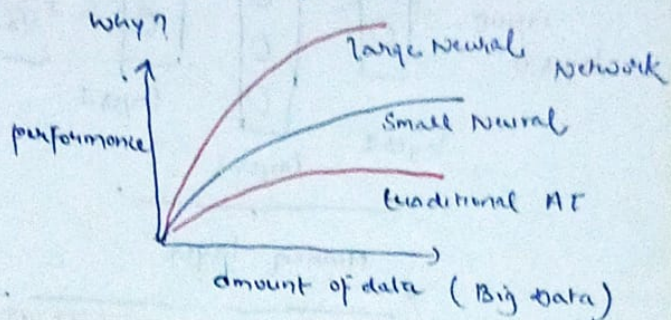
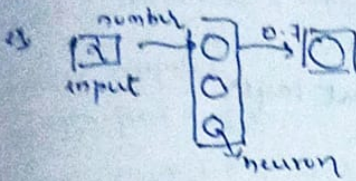


Week_2 Machine learning

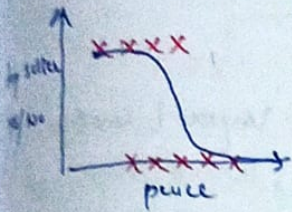
Advanced Learning Algorithms

Neural Networks

Algorithms that try to mimic the brain.



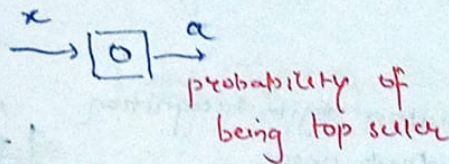
Demand Prediction



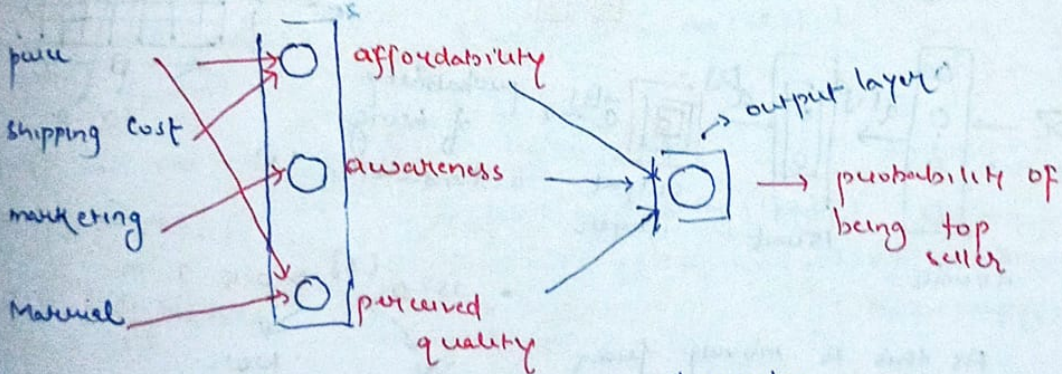
x = price input

$$a = f(x) = \frac{1}{1 + e^{-(wx + b)}}$$

activation output



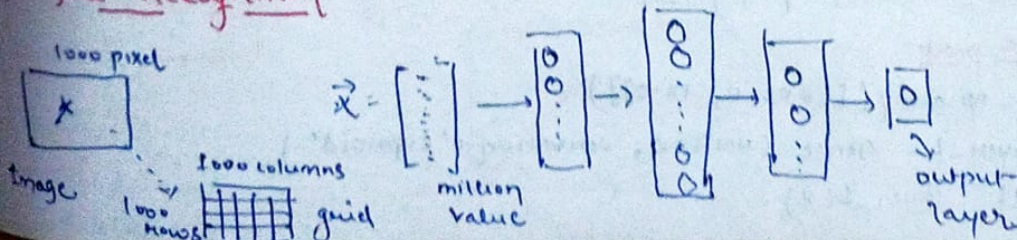
more features



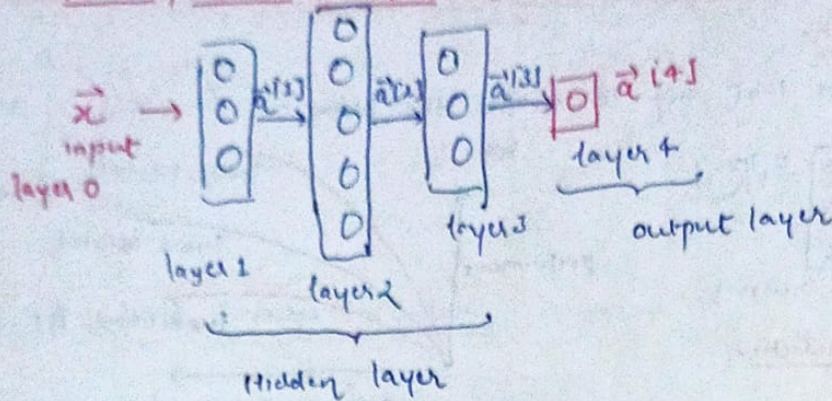
4 numbers
(input layer)

"activations"
3 numbers
(hidden layer)

Face Recognition



More complex neural network



Activation of value layer l

$$a_j^{[l]} = g(\bar{w}_j^{[l]} \cdot \bar{a}^{[l-1]} + b_j^{[l]})$$

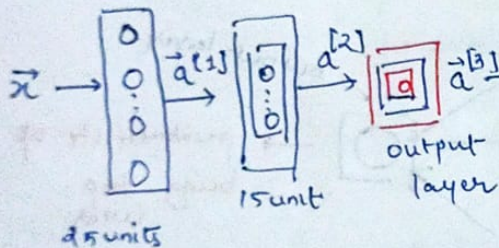
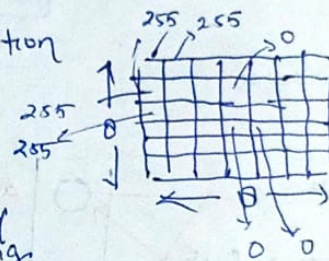
Parameter w & b of layer l unit j

{ sigmoid
"activation function" }

Forward Propagation

ex Handwritten digit recognition

Digit image



probability of being a handwritten '1'

• is $a_1^{[3]} \geq 0.5$?

yes $\hat{y} = 1$ no $\hat{y} = 0$

As this is moving from $a^{[1]} \rightarrow a^{[2]} \rightarrow a^{[3]}$

\therefore It is known as forward propagation.

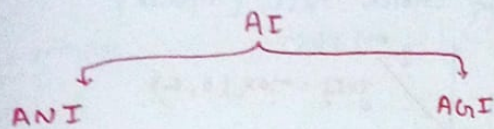
Code part:

```
x = np.array([1.0, 0.0, 1.0])
layer_1 = Dense(units=3, activation='sigmoid')
a1 = layer_1(x)
```


Building a neural network architecture

layer-1 = Dense (units=3, activation="sigmoid")
 layer-2 = Dense (units=1, activation="sigmoid")
 model = Sequential ([layer-1, layer-2])
 model.compile()
 model.fit(X, y)
 model.predict(X_new)

Syntax:



(Artificial narrow intelligence)

smart speaker, self driving car, web search, AI in farming and factories

(Artificial general intelligence)

do anything a human can do

2:

specify how to compute output given input x and parameters w, b (define model) $f_{w,b}(x) = ?$

specify loss and cost

$$L(f_{w,b}(\vec{x}), y)$$

$$J(w,b) = \frac{1}{m} \sum_{i=1}^m L(f_{w,b}(\vec{x}^{(i)}), y^{(i)})$$

Train on data to minimize cost function

Model Training Steps:
logistic regression

$$z = \text{np.dot}(w, x) + b$$

$$f_x = 1 / (1 + \text{np.exp}(-z))$$

logistic loss

$$\text{loss} = -y * \text{np.log}(f_x) - (1-y) * \text{np.log}(1-f_x)$$

$$w = w - \text{alpha} * dj - dw$$

$$b = b - \text{alpha} * dj - db$$

neural network

model = Sequential(
 [Dense(...)
 Dense(...)
 Dense(...)]

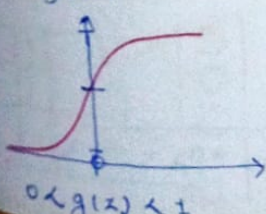
Binary cross Entropy

model.compile(
 loss=BinaryCrossentropy()

model.fit(X, y,
 epochs=100)

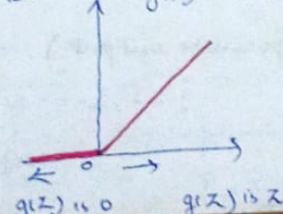
Alternatives to the Sigmoid Activation

Sigmoid

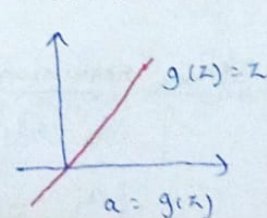


ReLU

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



Linear Activation

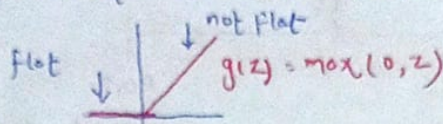


How to choose activation function for output layer or hidden layer

- Binary classification $\rightarrow \{0, 1\}$ problem \Rightarrow Sigmoid
 - Regression problem
e.g. tomorrow stock price \rightarrow Linear
 - e.g. House prediction but it will only have positive value \Rightarrow ReLU
- Output layer based on y prediction

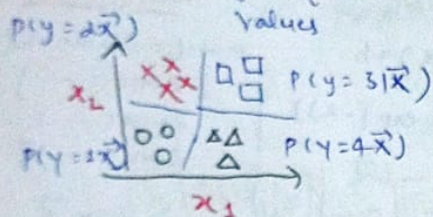
• Hidden Layer: most common choice ReLU {default}

- It is faster
- It has separate region of flat or not
- Not as smooth as sigmoid
- \therefore It learns pretty fast



• Multi class Classification

\downarrow target y can take more than just two 0 and 1 values



• Softmax:

Softmax regression (4 possible output)
 $y = 1, 2, 3, 4$

$$\otimes | z_1 = \vec{w}_1 \cdot \vec{x} + b_1 \quad a_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} = P(y=1|\vec{x})$$

$$\bigcirc | z_2 = \vec{w}_2 \cdot \vec{x} + b_2 \quad a_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} = P(y=2|\vec{x})$$

$$\square | z_3 = \vec{w}_3 \cdot \vec{x} + b_3 \quad = P(y=3|\vec{x})$$

$$\triangle | z_4 = \vec{w}_4 \cdot \vec{x} + b_4 \quad = P(y=4|\vec{x})$$

Softmax regression (N possible output)

$$z_j = \vec{w}_j \cdot \vec{x} + b_j \quad j = 1, 2, \dots, N$$

$$a_j = \frac{e^{z_j}}{\sum_{k=1}^N e^{z_k}}$$

Note: $a_1 + a_2 + \dots + a_N = 1$

Cost

Logistic regression

$$z = \vec{w} \cdot \vec{x} + b$$

$$a_1 = g(z) = \frac{1}{1 + e^{-z}} \Rightarrow P(y=1|\vec{x})$$

$$a_2 = 1 - a_1 = P(y=0|\vec{x})$$

$$loss = -y \log a_1 - (1-y) \log (1-a_1)$$

$J(\vec{w}, b)$ = average loss

python $f_x = \text{tf.nn.sigmoid}(\text{logit})$
 something similar like this

Softmax regression

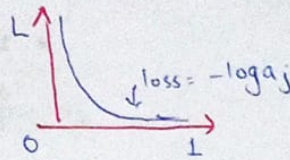
$$a_1 = \frac{e^{z_1}}{a_1 e^{z_1} + e^{z_2} + \dots + e^{z_n}}$$

$$\vdots$$

$$a_n$$

Crossentropy loss

$$loss(a_1, \dots, a_n) = \begin{cases} -\log a_1 & \text{if } y=1 \\ -\log a_2 & \text{if } y=2 \\ \vdots \end{cases}$$



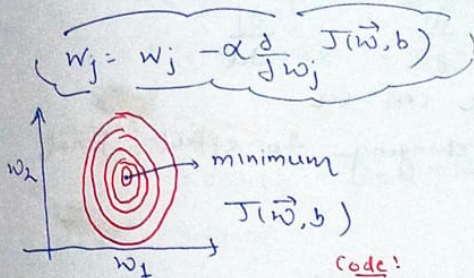
* model.compile(loss = Sparse Categorical Cross entropy())

Rest is same except change in outer layer.

Adam Algorithm

It increases the learning rate α when it sees that gradient descent taking very minute steps to reach minima

• It also decreases alpha when learning rate is too fast



Adam: Adaptive Moment estimation

It adjust α automatically.

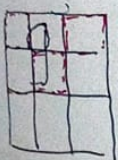
model.compile(optimizer = tf.keras.optimizers.Adam (

learning rate = $1e-3$)

$\Rightarrow \alpha = 10^{-3}$

Additional Layer Types

Convolutional Layer



Each neuron only look at part of the previous layer's output { each square in grid }

→ faster computation

→ Need less training data.

How Back Propagation Works {calculus part}

• Derivative of

Cost function $J(w) = w^2$

say $w = 3$ $J(w) = 3^2 = 9$

If we increase w by a tiny amount $\epsilon = 0.001$ how does $J(w)$ change?

$$w = 3 + 0.001$$

$$J(w) = w^2 = 9.006001$$

If $w \uparrow 0.001$

$$J(w) \uparrow 6 \times 0.001 \quad 6 \times \epsilon$$

$$\frac{\partial}{\partial w} J(w) = 6$$

* If $w \uparrow \epsilon$ cause $J(w) \uparrow k \times \epsilon$ then

$$\frac{\partial}{\partial w} J(w) = k$$

gradient descent

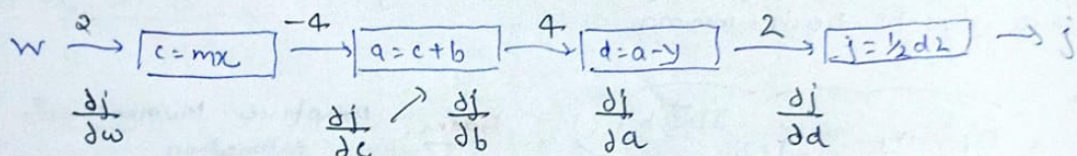
$$\text{repeat } \{ \quad w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b) \}$$

Q. why we use backpropagation to calculate derivatives

If derivative is small, then this update step will make a small update on w_j

If large \rightarrow so large update

* Backprop is an efficient way to compute derivatives



By calculating right to left we can see

how these quantities are changing to effect final output j