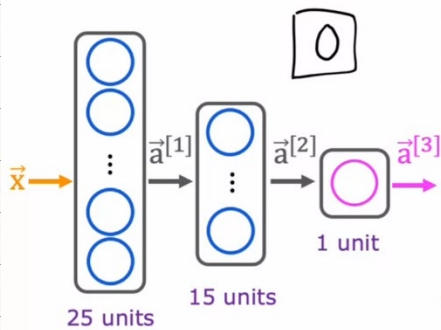


Subject: / /

Train a Neural Network:

Train a Neural Network in TensorFlow



```
import tensorflow as tf
from tensorflow.keras import Sequential
from tensorflow.keras.layers import Dense

model = Sequential([
    Dense(units=25, activation='sigmoid'),
    Dense(units=15, activation='sigmoid'),
    Dense(units=1, activation='sigmoid'),
])
from tensorflow.keras.losses import BinaryCrossentropy
model.compile(loss=BinaryCrossentropy())
model.fit(X, Y, epochs=100)
```

①
②
③
epochs: number of steps in gradient descent

Given set of (x, y) examples

How to build and train this in code?

Training Details:

① Specify how to compute output given input x and parameters w, b
 $f_{w,b}(\vec{x}) = ?$

② Specify loss and Cost

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

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

③ Train on data to minimize $J(\vec{w}, b)$

For a neural network,

① `model = Sequential(...)`

② `model.compile(loss=BinaryCrossentropy())`

③ `model.fit(X, y, epochs=100)`

Subject: / /

Loss and Cost function:

Classification:

$$L(f(\vec{x}), y) = -y \log(f(\vec{x})) - (1-y) \log(1-f(\vec{x})) \rightarrow \text{Binary cross entropy}$$

model.compile(loss = BinaryCrossentropy())

Regression:

model.compile(loss = MeanSquaredError)

$$J(w, b) = \frac{1}{n} \sum_{i=0}^n L(f(\vec{x}^{(i)}), y^{(i)})$$

$w^{(0)}, w^{(1)}, w^{(2)}$ $b^{(0)}, b^{(1)}, b^{(2)}$

3. Gradient descent:

Repeat {

$$w_j^{(k+1)} = w_j^{(k)} - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)$$

$$b_j^{(k+1)} = b_j^{(k)} - \alpha \frac{\partial}{\partial b_j} J(\vec{w}, b)$$

Backpropagation
to compute the
derivatives

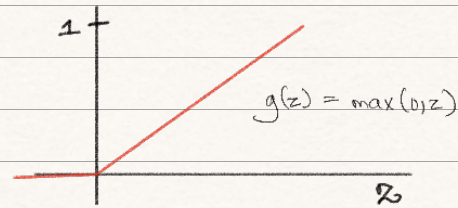
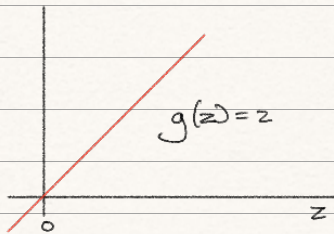
}

Alternatives to Sigmoid activation:

Subject:

Linear activation function

ReLU - Rectified Linear Unit



How to choose?

- Depends on target
- For Binary classification, use sigmoid
- For Regression, use Linear activation function if y is $+/-$.
use ReLU if y is $0/+$.
- For hidden layers, use ReLU

Multi-class classification:

↳ More than one class/category

Soft-max:

↳ Generalization of logistic Regression

$$X \quad z_1 = \vec{w}_1 \cdot \vec{x} + b_1$$

$$O \quad z_2 = \vec{w}_2 \cdot \vec{x} + b_2$$

$$\square \quad z_3 = \vec{w}_3 \cdot \vec{x} + b_3$$

$$\Delta \quad z_4 = \vec{w}_4 \cdot \vec{x} + b_4$$

$$p(y=1|\vec{x}) = a_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

\vdots

$$p(y=4|\vec{x}) = a_4 = \frac{e^{z_4}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

Subject:

//

$$y = 1, 2, 3, \dots, N$$

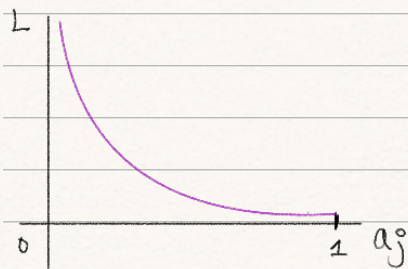
$$z_j = \vec{w}_j \cdot \vec{x} + b_j \quad j = 1, \dots, N$$

$$a_j = \frac{e^{z_j}}{\sum_{k=1}^N e^{z_k}} = P(y=j | \vec{x})$$

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

Cost:

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



- Loss is more if a_j is small
- Motivates model to make a_j as close as possible to 1.

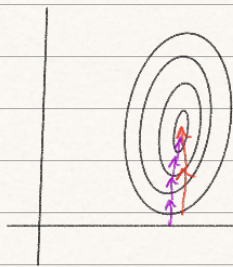
$$\text{loss} = \text{SparseCategoricalCrossentropy}(l)$$

Subject:

Multiclass is classifying single piece, y is scalar
Multi-label is classifying many pieces, y is vector

$$J(w, b) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{j=1}^N 1\{y^{(i)} = j\} \log \frac{e^{z_j^{(i)}}}{\sum_{k=1}^N e^{z_k^{(i)}}} \right]$$

Advanced propagation:



α is big } done automatically
 α is small

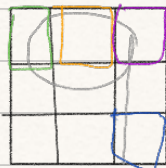
as gradient descent is
in same direction, why not
make α bigger? Adam does
this.

Adam:

- Adaptive Moment estimation
- α is diff for each parameter

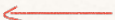
Additional layer types:

Convolutional layer

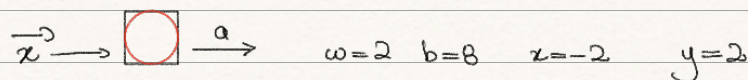


Subject: / /



- Helps avoid overfitting

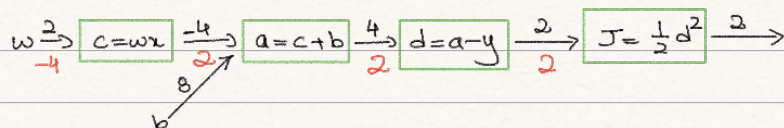
Back propagation:  right - left

Computation graph:



$$J(w, b) = \frac{1}{2} (a - y)^2$$

 Forward
 back
prop



$$\frac{\partial J}{\partial d} = 2$$

$$\frac{\partial J}{\partial a} = 2$$

$$d \uparrow 0.001 \quad J \uparrow 0.002$$

$$a \uparrow 0.001 \quad d \uparrow 0.001 \quad J \uparrow 0.002$$

$$a = 4.001 \quad d = 2.001$$

$$\frac{\partial J}{\partial a} = \frac{\partial d}{\partial a} \times \frac{\partial J}{\partial d}$$

$$= 1 \times 2$$

$$\frac{\partial J}{\partial c} = 2$$

$$\frac{\partial J}{\partial b} = 2$$

$$\frac{\partial J}{\partial w} = -4$$

$$c \uparrow 0.001 \quad a \uparrow 0.001$$

$$J \uparrow 0.002$$

$$b \uparrow 0.001 \quad a \uparrow 0.001$$

$$J \uparrow 0.002$$

$$w \uparrow 0.001 \quad c = -4.002$$

$$w = 2.001 \quad c \downarrow 2 \times 0.001$$

$$c \uparrow -2 \times 0.001$$

$$\frac{\partial J}{\partial c} = \frac{\partial a}{\partial c} \times \frac{\partial J}{\partial a}$$

$$= 1 \times 2$$

$$\frac{\partial J}{\partial b} = \frac{\partial a}{\partial b} \times \frac{\partial J}{\partial a}$$

$$= 1 \times 2$$

$$J \uparrow -4 \times 0.001$$

$$\frac{\partial J}{\partial w} = \frac{\partial c}{\partial w} \times \frac{\partial J}{\partial c}$$

$$= -2 \times 2$$

Subject: _____

/ /

$$w \xrightarrow{\exists} a = 2 + 3w \xrightarrow{\parallel} J = a^2 \xrightarrow{121} J$$

$$\frac{2w}{2J} = 66$$

$$\frac{2J}{2a} = 22$$

$$a \uparrow 0.001 \quad J \uparrow 0.022 \Rightarrow 22 \times 0.001$$

$$w \uparrow 0.001 \quad a \uparrow 0.003 \quad a = 11.001 \quad J = 121.022$$

$$3 \times 0.001$$