

Day-29, Oct-31, 2024.

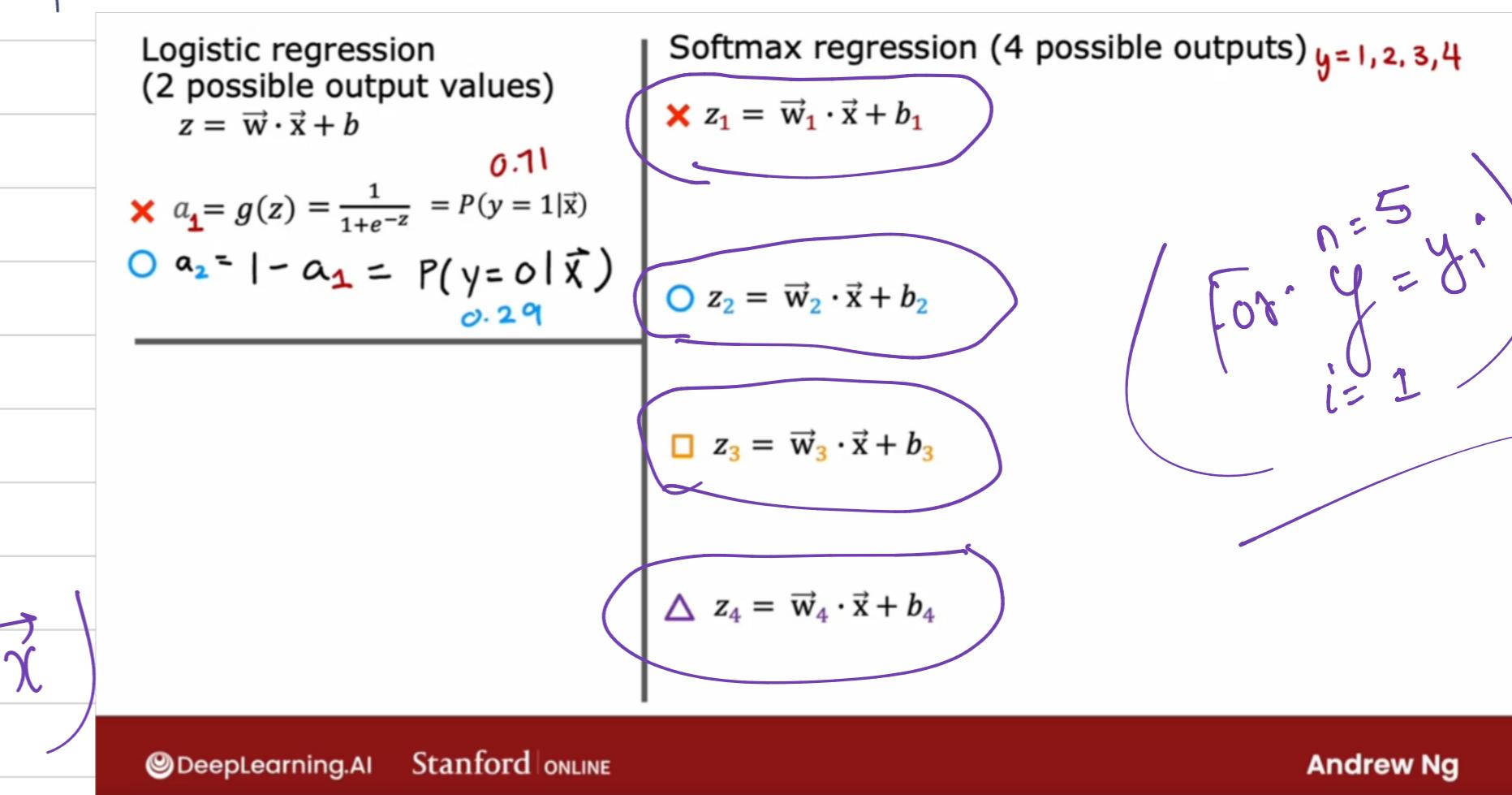
Softmax Activation Functions.

we have only 2 possible output values in logistic regression

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

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

Suppose $y=1$ given then if a_1 is 0.7 then



$$\text{if } a_1 \text{ is } 0.7 \text{ then } a_2 = 1 - a_1 \Rightarrow P(y=0 | \vec{x}) = \frac{(0 \cdot 3)}{1.0}$$

$$\begin{aligned} a_1 &= 0.7 \\ a_2 &= 0.3 \\ P &\rightarrow 1.0 \end{aligned}$$

So for $y = 1, 2, 3, 4$

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

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

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

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

So a_1, a_2, a_3 & a_4 are

$$a_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

$$a_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} \Rightarrow P(y=2|\vec{x})$$

Logistic regression
(2 possible output values)

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

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

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

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

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

$$a_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

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

$$a_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

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

$$a_3 = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

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

$$a_4 = \frac{e^{z_4}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

DeepLearning.AI

Stanford ONLINE

Andrew Ng

$$\cdots \left\{ \begin{array}{l} a_4 = \frac{e^{z_4}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} \\ \Rightarrow P(y=4|\vec{x}) \end{array} \right.$$

Logistic regression
(2 possible output values)

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

$$\times z_1 = \vec{w}_1 \cdot \vec{x} + b_1 \quad 0.11$$

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

$$\circ a_2 = 1 - a_1 = P(y=0|\vec{x}) \quad 0.29$$

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

$$\times z_1 = \vec{w}_1 \cdot \vec{x} + b_1 \quad a_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

$$= P(y=1|\vec{x}) \quad 0.30$$

$$\circ z_2 = \vec{w}_2 \cdot \vec{x} + b_2 \quad a_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

$$= P(y=2|\vec{x}) \quad 0.20$$

$$\square z_3 = \vec{w}_3 \cdot \vec{x} + b_3 \quad a_3 = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

$$= P(y=3|\vec{x}) \quad 0.15$$

$$\Delta z_4 = \vec{w}_4 \cdot \vec{x} + b_4 \quad a_4 = \frac{e^{z_4}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

$$= P(y=4|\vec{x}) \quad 0.35$$

Logistic regression
(2 possible output values)

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

$$\times z_1 = \vec{w}_1 \cdot \vec{x} + b_1 \quad 0.11$$

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

$$\circ a_2 = 1 - a_1 = P(y=0|\vec{x}) \quad 0.29$$

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

$$\times z_1 = \vec{w}_1 \cdot \vec{x} + b_1 \quad a_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

$$= P(y=1|\vec{x}) \quad 0.30$$

$$\circ z_2 = \vec{w}_2 \cdot \vec{x} + b_2 \quad a_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

$$= P(y=2|\vec{x}) \quad 0.20$$

$$\square z_3 = \vec{w}_3 \cdot \vec{x} + b_3 \quad a_3 = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

$$= P(y=3|\vec{x}) \quad 0.15$$

$$\Delta z_4 = \vec{w}_4 \cdot \vec{x} + b_4 \quad a_4 = \frac{e^{z_4}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

$$= P(y=4|\vec{x}) \quad 0.35$$

\sum of a_j must be equal to 1 because probability is always $\frac{1}{1}$.

Notation for Activation Function (Soft Max)

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

$$\Rightarrow P(y=j|\vec{x})$$

Recall.

$$a_1 = g(z) \Rightarrow \frac{1}{1 + e^{-z}}$$

$$\Rightarrow P(y=1 | \vec{x})$$

$$a_2 = 1 - a_1$$

$$\Rightarrow P(y=0 | \vec{x})$$

$$\text{loss} = -y \log a_1 - (1-y) \log (1-a_1)$$

if $y=1$

if $y=0$

for SoftMax Regression? So, we figure out & observe the Crossentropy loss (predicted distribution probability to the original).

Logistic regression

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

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

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

$$\text{loss} = -y \log a_1 - (1-y) \log (1-a_1)$$

if $y=1$

Cost

Softmax regression

DeepLearning.AI Stanford ONLINE

Andrew Ng

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

Softmax Regression

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

$$\Rightarrow P(y=1|\vec{x})$$

$$a_N \Rightarrow \frac{e^{z_N}}{e^{z_1} + e^{z_2} + \dots + e^{z_N}}$$

$$\Rightarrow P(y=N|\vec{x})$$

Logistic regression

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

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

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

$$\text{loss} = -y \underbrace{\log a_1}_{\text{if } y=1} - (1-y) \underbrace{\log(1-a_1)}_{\text{if } y=0}$$

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

Cost

Softmax regression

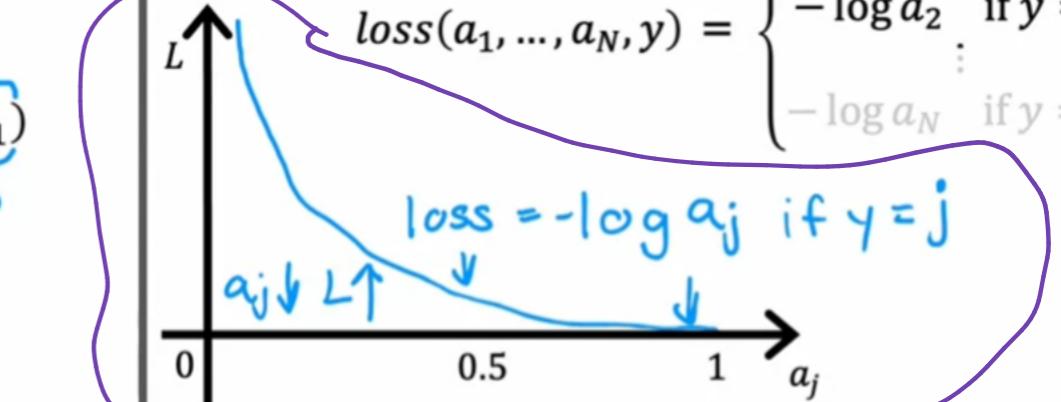
$$a_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + \dots + e^{z_N}} = P(y=1|\vec{x})$$

$$\vdots$$

$$a_N = \frac{e^{z_N}}{e^{z_1} + e^{z_2} + \dots + e^{z_N}} = P(y=N|\vec{x})$$

Crossentropy loss

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



As you go near the '1' a_j if is loss = $-\log a_j$
 but if go near the '0' it is $\log(1-a_j)$

Neural Network with Softmax output

ReLU used in 25 nodes

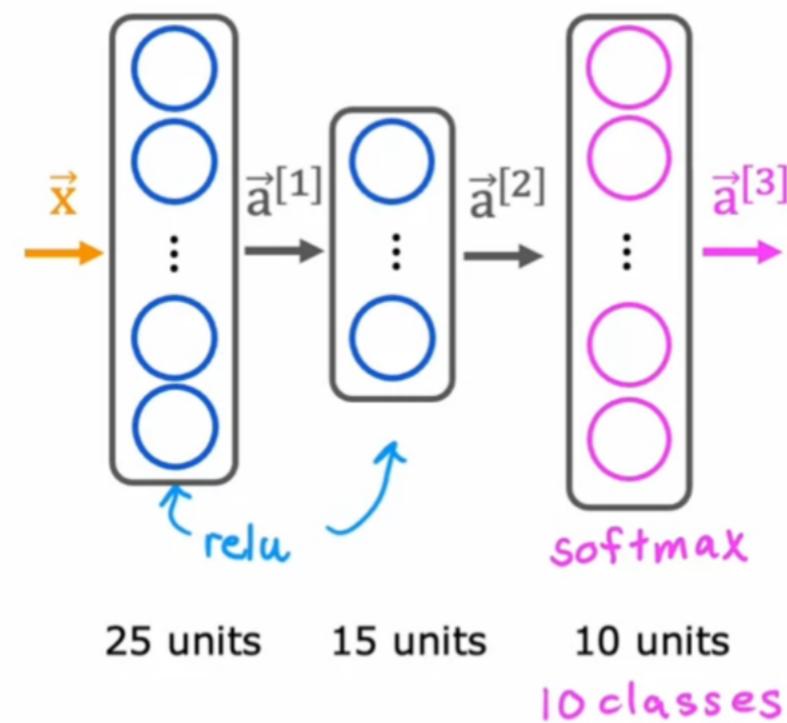
and 15 nodes and

they are hidden layers.

Whereas Softmax is

used in output layer

having 10 nodes or



DeepLearning.AI Stanford ONLINE

Andrew Ng

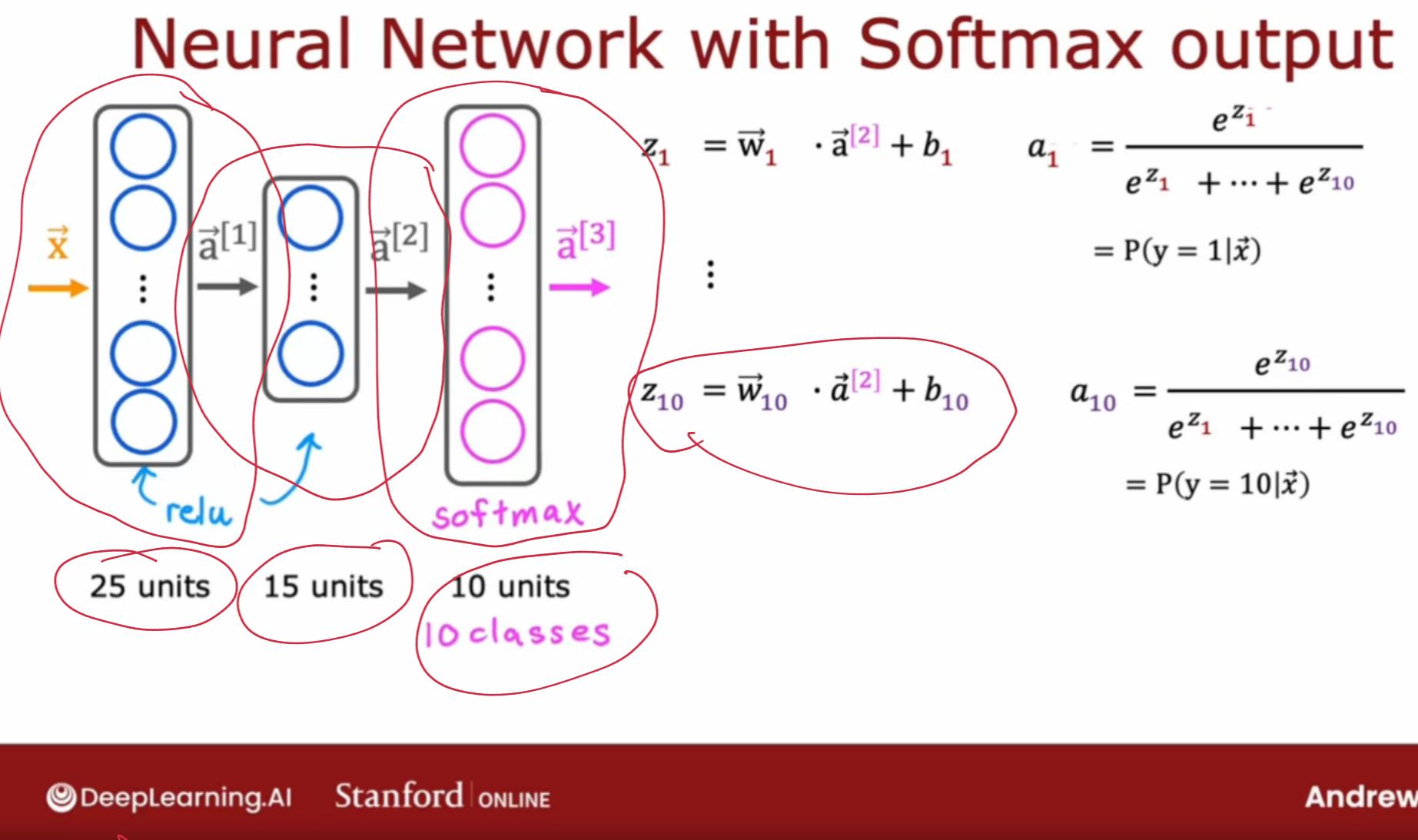
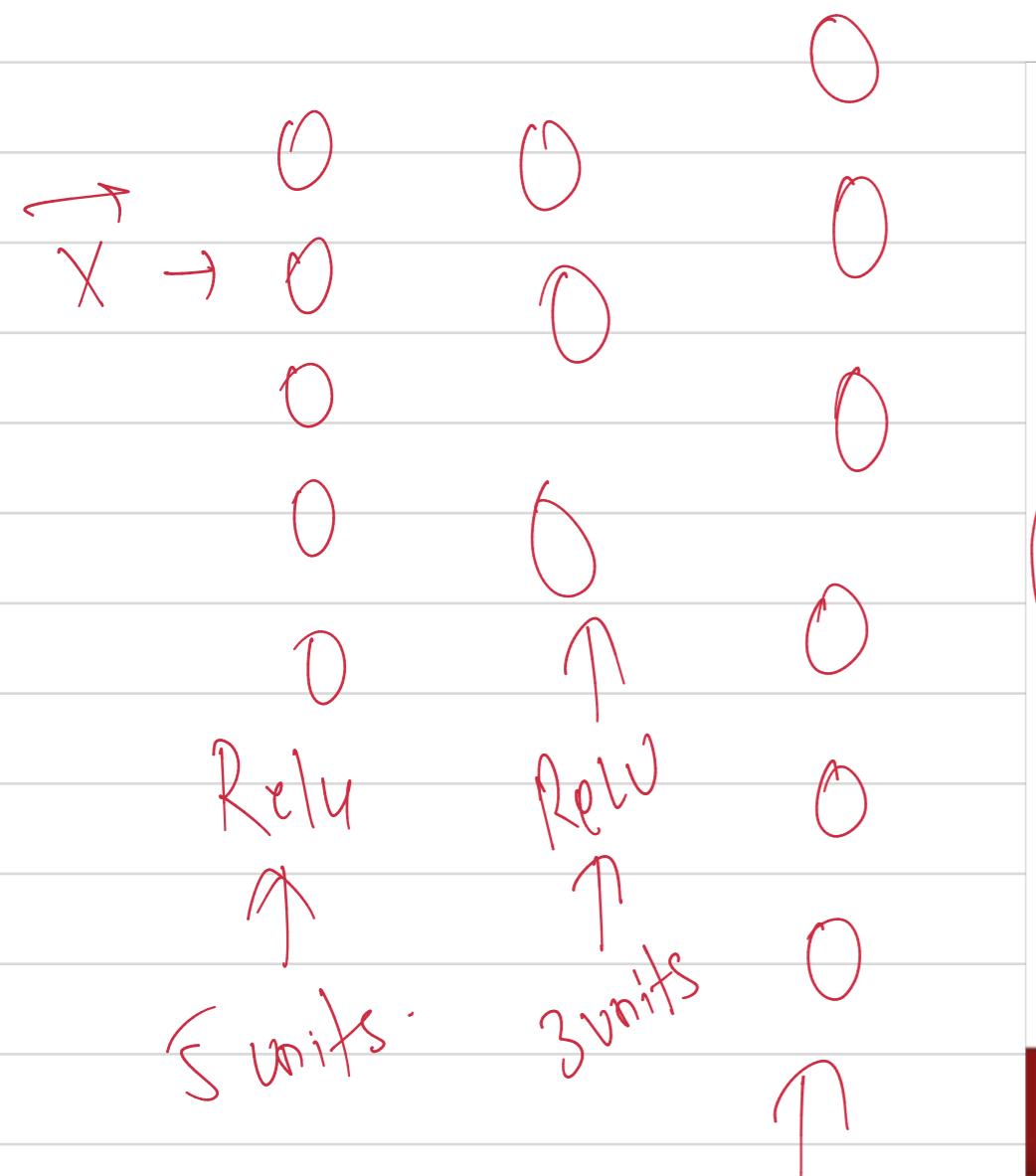
10 classes →

for $z_n = 5$

Output class
5

$$z_1 = \vec{w}_1 \cdot \vec{a}^{[2]} + b_1 \Rightarrow a_1 \Rightarrow \frac{e^{z_1}}{e^{z_1} + \dots + e^{z_{10}}} = p(y=1 | \vec{x})$$

$$z_5: \vec{w}_5 \cdot \vec{a}^{[2]} + b_5 \Rightarrow a_5 \Rightarrow \frac{e^{z_5}}{e^{z_1} + \dots + e^{z_{10}}} = p(y=5 | \vec{x})$$



So, we add Superscript to denote the layer of Neural Network
Whereas Subscript tells the index in the particular layer itself.

MNIST with softmax

① specify the model

$$f_{\vec{w}, \vec{b}}(\vec{x}) = ?$$

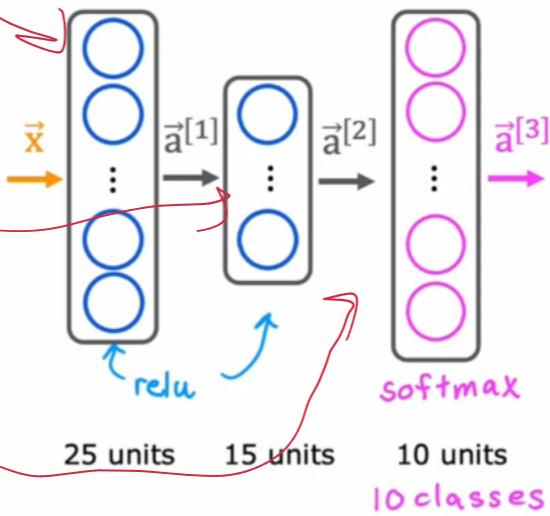
```
import tensorflow as tf
from tensorflow.keras import Sequential
from tensorflow.keras.layers import Dense
model = Sequential([
    Dense(units=25, activation='relu'),
    Dense(units=15, activation='relu'),
    Dense(units=10, activation='softmax')
])
```

specify loss and cost

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

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

Neural Network with Softmax output



$$z_1^{[3]} = \vec{w}_1^{[3]} \cdot \vec{a}^{[2]} + b_1^{[3]} \quad a_1^{[3]} = \frac{e^{z_1^{[3]}}}{e^{z_1^{[3]}} + \dots + e^{z_{10}^{[3]}}}$$

$$= P(y = 1 | \vec{x})$$

$$z_{10}^{[3]} = \vec{w}_{10}^{[3]} \cdot \vec{a}^{[2]} + b_{10}^{[3]} \quad a_{10}^{[3]} = \frac{e^{z_{10}^{[3]}}}{e^{z_1^{[3]}} + \dots + e^{z_{10}^{[3]}}}$$

$$= P(y = 10 | \vec{x})$$

logistic regression

$$a_1^{[3]} = g(z_1^{[3]}) \quad a_2^{[3]} = g(z_2^{[3]})$$

softmax

$$\vec{a}^{[3]} = (a_1^{[3]}, \dots, a_{10}^{[3]}) = g(z_1^{[3]}, \dots, z_{10}^{[3]})$$

Loss function
Sparse Categorical Cross Entropy
 $\log(y)$ used in
classification - dealing with
multi-class classification problem
Input \rightarrow 1
Category $[0, 10] \rightarrow$ Sparse
only select or predict category 1.
takes [0 to 10]
&
only sparse one value it
can be only one of the
category for a given digit

MNIST with softmax

① specify the model

$$f_{\vec{w}, b}(\vec{x}) = ?$$

② specify loss and cost

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

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

```
import tensorflow as tf
from tensorflow.keras import Sequential
from tensorflow.keras.layers import Dense
model = Sequential([
    Dense(units=25, activation='relu'),
    Dense(units=15, activation='relu'),
    Dense(units=10, activation='softmax')
])
```

```
from tensorflow.keras.losses import
SparseCategoricalCrossentropy
```

```
model.compile(loss= SparseCategoricalCrossentropy() )
```

① specify the model

$$f_{\vec{w}, b}(\vec{x}) = ?$$

② specify loss and cost

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

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

MNIST with softmax

```
import tensorflow as tf
from tensorflow.keras import Sequential
from tensorflow.keras.layers import Dense
model = Sequential([
```

```
    Dense(units=25, activation='relu'),
    Dense(units=15, activation='relu'),
    Dense(units=10, activation='softmax')
])
```

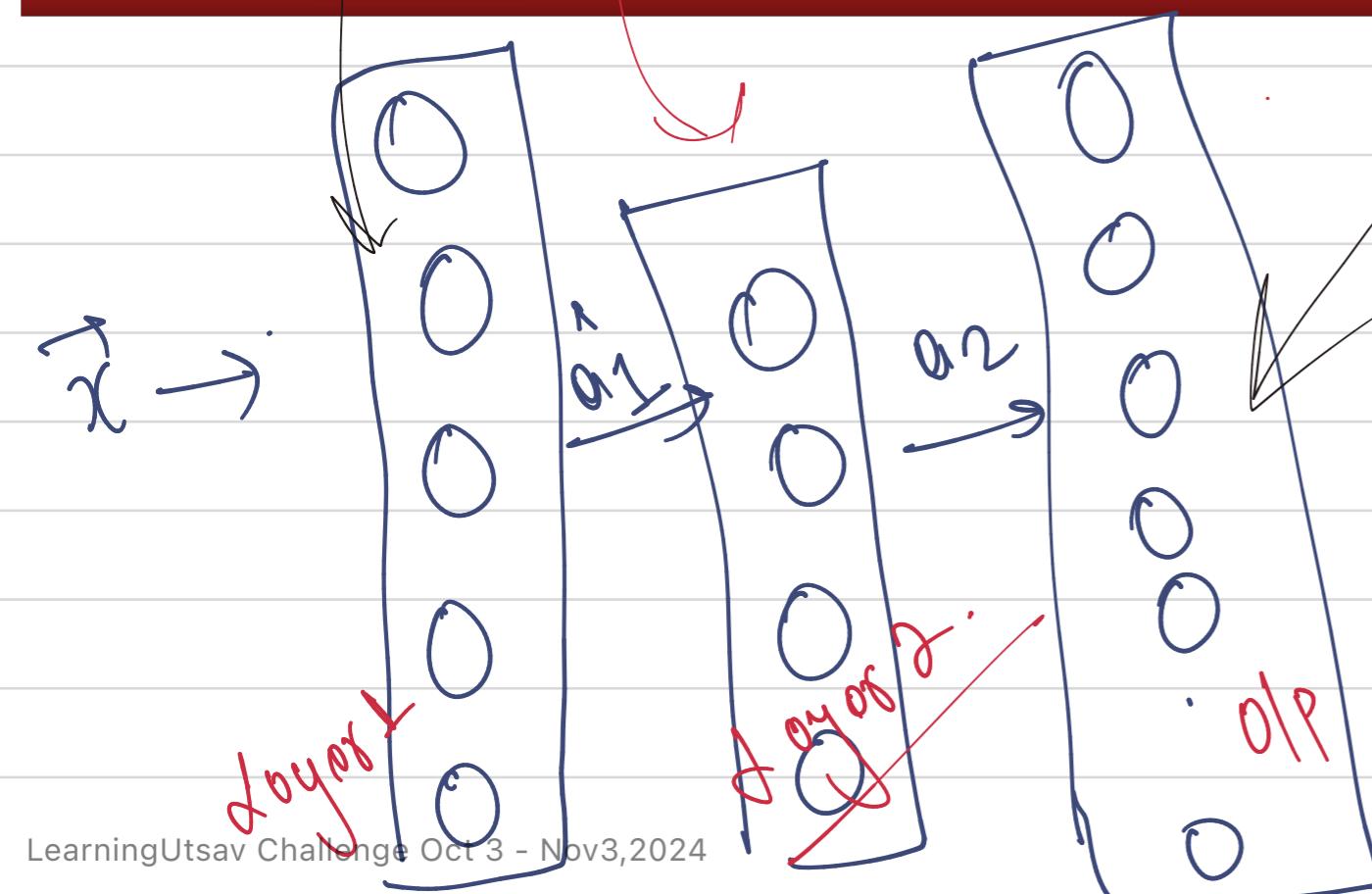
```
from tensorflow.keras.losses import
SparseCategoricalCrossentropy
```

```
model.compile(loss= SparseCategoricalCrossentropy() )
```

```
model.fit(X, Y, epochs=100)
```

Note: better (recommended) version later.

Don't use the version shown here!



So, the main idea from the slides
is that we use softmax for
multi-class supported by
Sparse Categorical CrossEntropy (loss
function).