

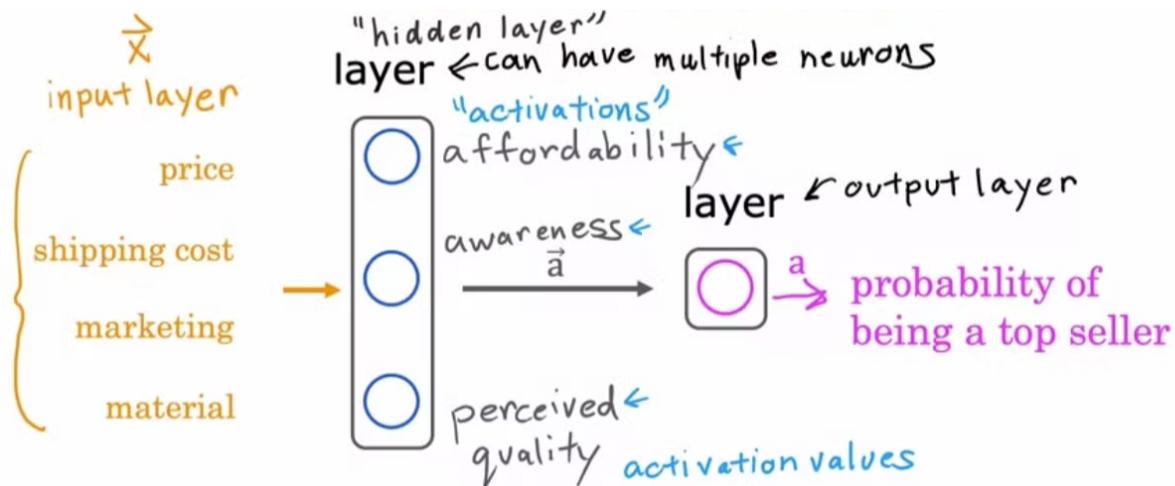
Advanced Learning Algorithms

Origins: algorithms that try to mimic the brain

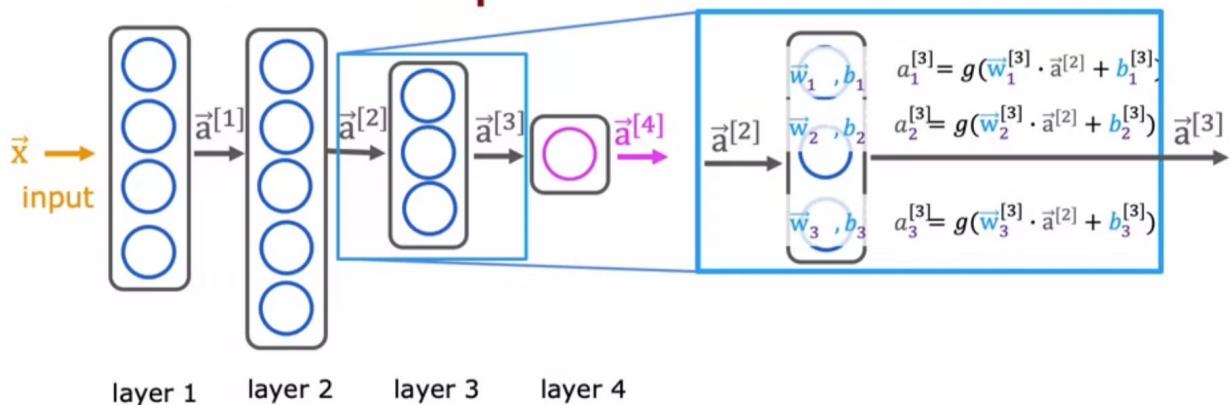
Ressurge from around 2005 under the term *deep learning*, as the amount of data rapidly increased. The rise of GPUs was also a major force.

Applications: speech recognition → computer vision → NLP ...

Example: predict if t-shirt will become a top-seller



It's like a version of logistic regression that learns its own features layer by layer. That replaces manual feature engineering.

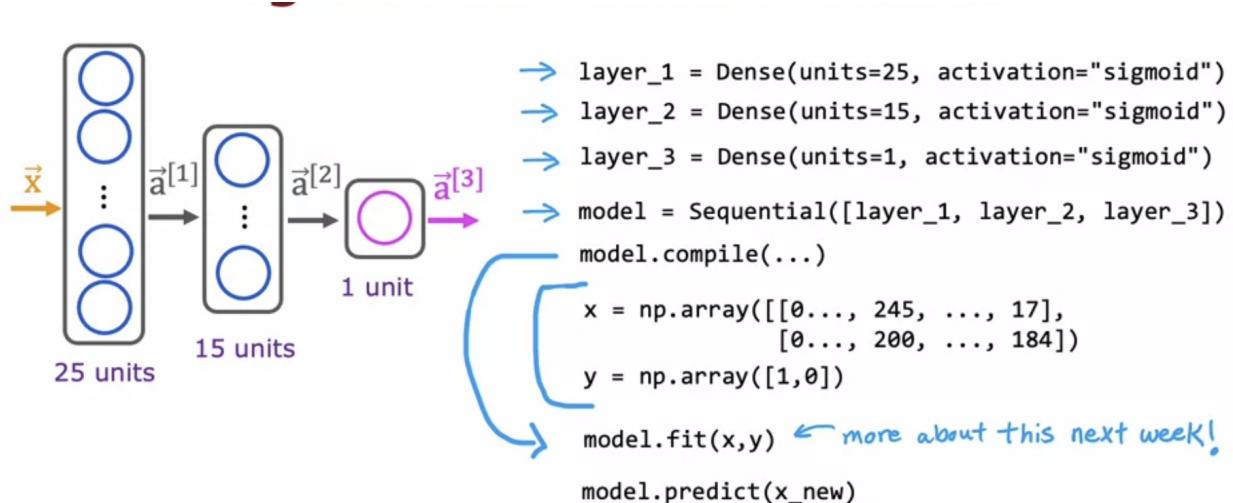
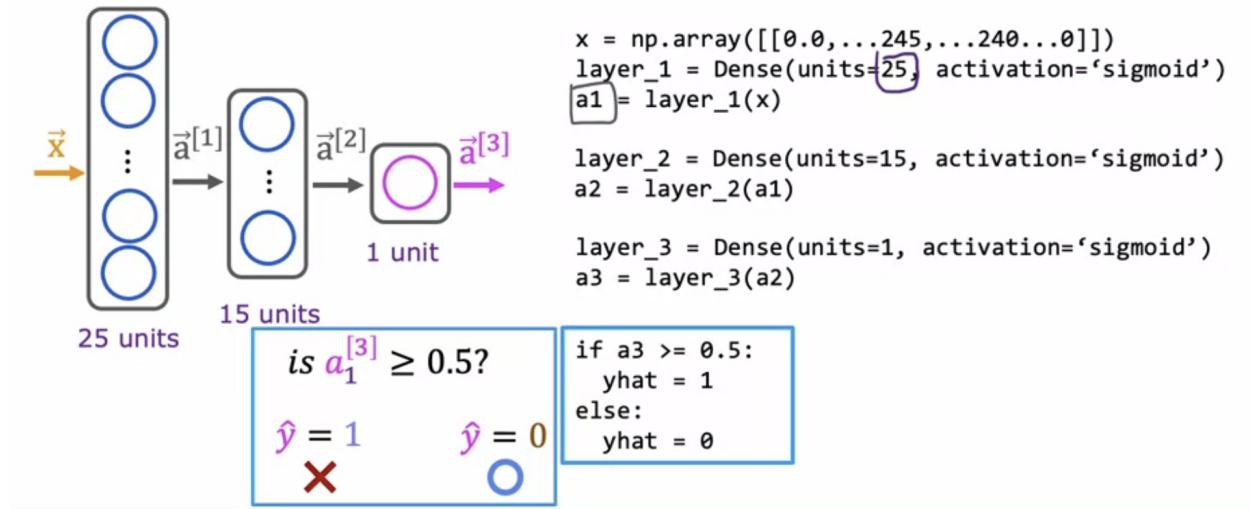


The activation function g could be the sigmoid function or other functions.

Forward propagation to make a prediction (inference).

TensorFlow

Model for digit classification



Vectorized implementation makes neural networks run much faster!

Dense layer vectorized

$A^T = [200 \quad 17]$

$W = \begin{bmatrix} 1 & -3 & 5 \\ -2 & 4 & -6 \end{bmatrix}$

$B = [-1 \quad 1 \quad 2]$

$Z = A^T W + B$

$A = g(Z)$

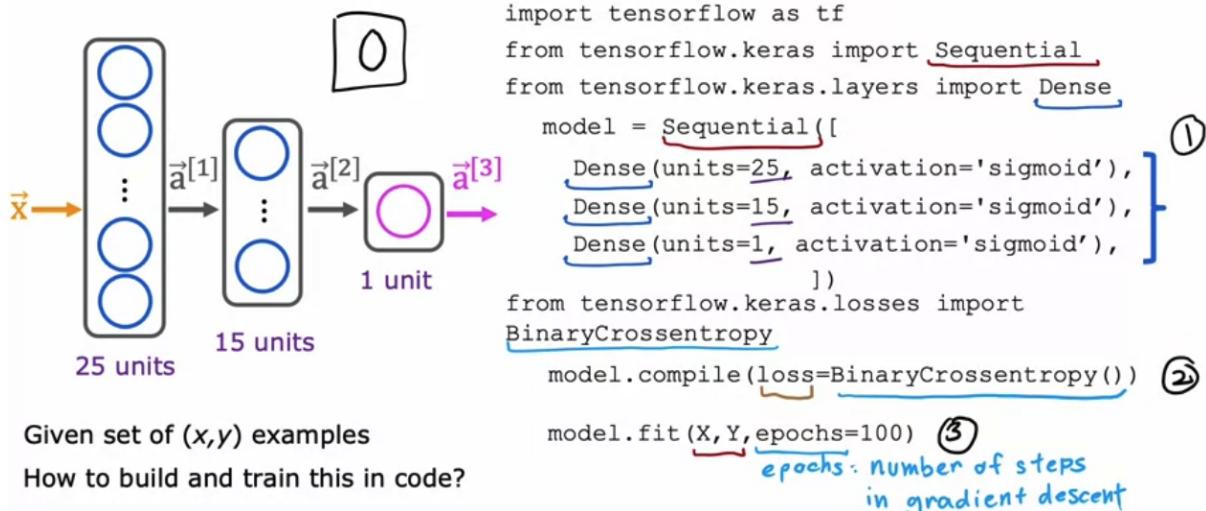
$A^2 = \begin{bmatrix} 165 & -531 & 900 \end{bmatrix}$

```

A
AT = np.array([[200, 17]])
W = np.array([[1, -3, 5],
[-2, 4, -6]])
b = np.array([[-1, 1, 2]])
a_in
def dense(AT,W,b):
    z = np.matmul(AT,W) + b
    a_out = g(z)
    return a_out
[[1,0,1]]

```

Train a Neural Network in TensorFlow



Model Training Steps

TensorFlow

<p>① specify how to compute output given input x and parameters w, b (define model) $f_{\bar{w}, \bar{b}}(\vec{x}) = ?$</p> <p>② specify loss and cost $L(f_{\bar{w}, \bar{b}}(\vec{x}), y)$ 1 example $J(\bar{w}, \bar{b}) = \frac{1}{m} \sum_{i=1}^m L(f_{\bar{w}, \bar{b}}(\vec{x}^{(i)}), y^{(i)})$</p> <p>③ Train on data to minimize $J(\bar{w}, \bar{b})$</p>	<p>logistic regression</p> <pre><code>z = np.dot(w, x) + b f_x = 1 / (1 + np.exp(-z))</code></pre> <p>logistic loss</p> <pre><code>loss = -y * np.log(f_x) - (1 - y) * np.log(1 - f_x)</code></pre> <p>w = w - alpha * dj_dw b = b - alpha * dj_db</p>	<p>neural network</p> <pre><code>model = Sequential([Dense(...), Dense(...), Dense(...),])</code></pre> <p>binary cross entropy</p> <pre><code>model.compile(loss=BinaryCrossentropy())</code></pre> <p>model.fit(X, y, epochs=100)</p>
---	---	--

2. Loss and cost functions

handwritten digit classification problem binary classification

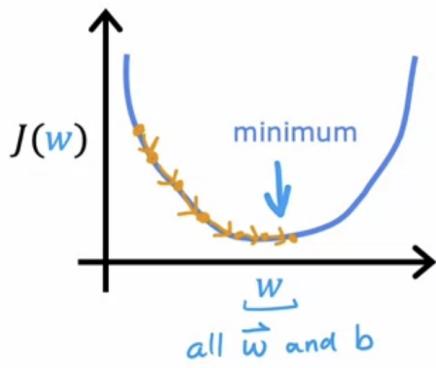
$L(f(\vec{x}), y) = -y \log(f(\vec{x})) - (1 - y) \log(1 - f(\vec{x}))$

Compare prediction vs. target

logistic loss
also known as binary cross entropy

model.compile(loss= BinaryCrossentropy()) from tensorflow.keras.losses import
BinaryCrossentropy Keras

3. Gradient descent



```

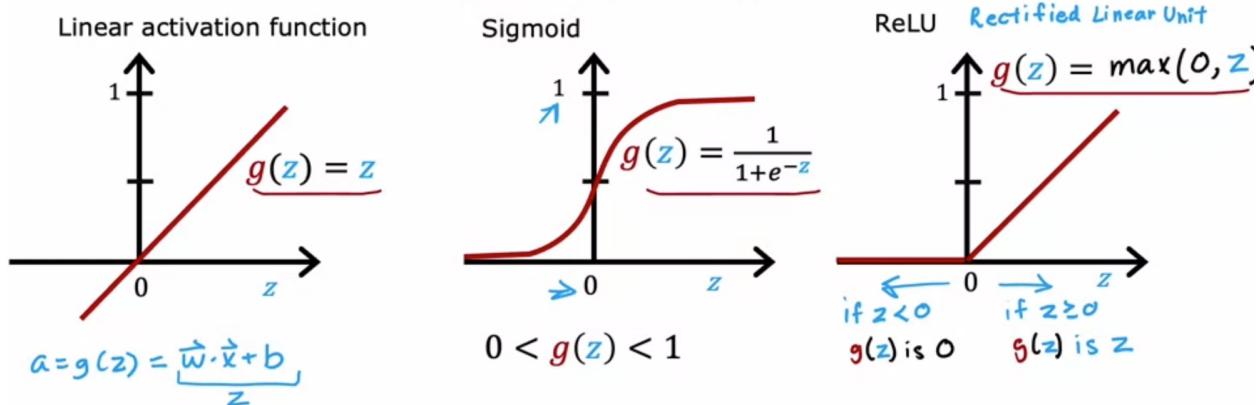
repeat {
     $w_j^{[l]} = w_j^{[l]} - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)$ 
     $b_j^{[l]} = b_j^{[l]} - \alpha \frac{\partial}{\partial b_j} J(\vec{w}, b)$ 
}
} Compute derivatives
for gradient descent
using "back propagation"

```

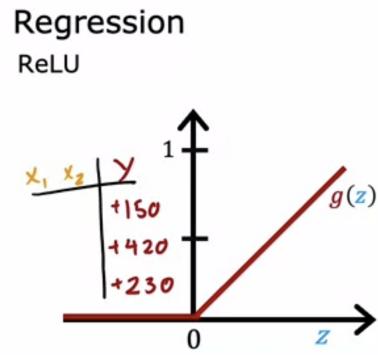
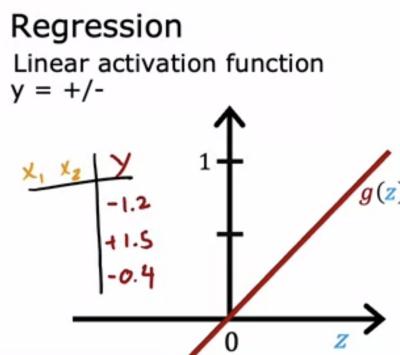
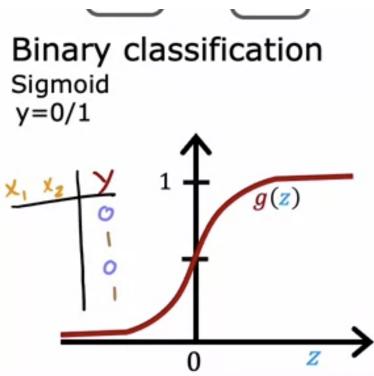
`model.fit(X, y, epochs=100)`

Another popular activation function is **ReLU** (rectified linear unit).

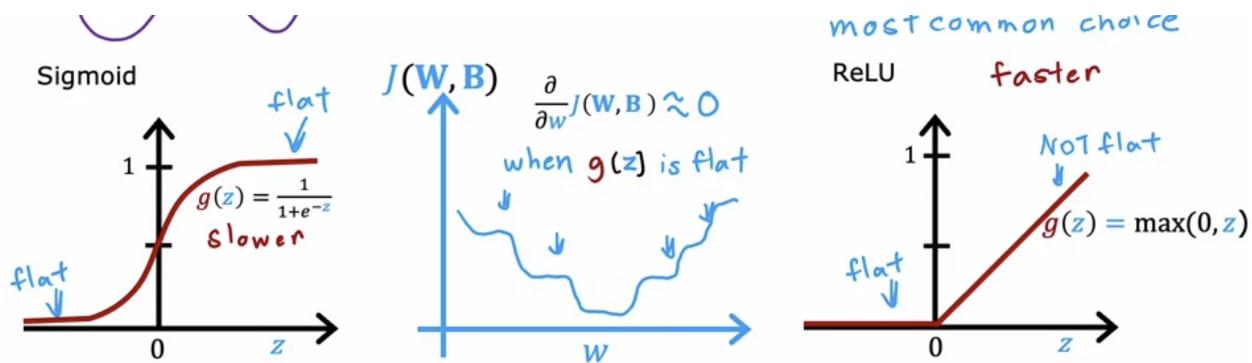
$$g(z) = \max(0, z)$$



We choose the activation function for the output layer depending on the nature of the label y we're trying to predict.



For the hidden layers, ReLU is the most commonly used. It's faster to compute and makes gradient descent faster. This is because ReLU is flat to the left only, whereas sigmoid is flat to the left and right.



Why don't we use linear activation function in the entire neural network? Because that would be equivalent to a single linear regression. And if the hidden layers are all linear and the output layer is sigmoid, then the neural network is equivalent to logistic regression.

The "off" or disable feature of the ReLU activation enables models to stitch together linear segments to model complex non-linear functions.

Multiclass classification

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

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

$\text{X} \quad \text{O} \quad \square \quad \Delta$
 $= P(y = 1|\vec{x})$

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

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

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

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

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

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

Softmax regression
(N possible outputs) $y=1, 2, 3, \dots, N$

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

parameters w_1, w_2, \dots, w_N
 $e^{z_j} \quad b_1, b_2, \dots, b_N$

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

Cost

Logistic regression

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

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

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

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

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

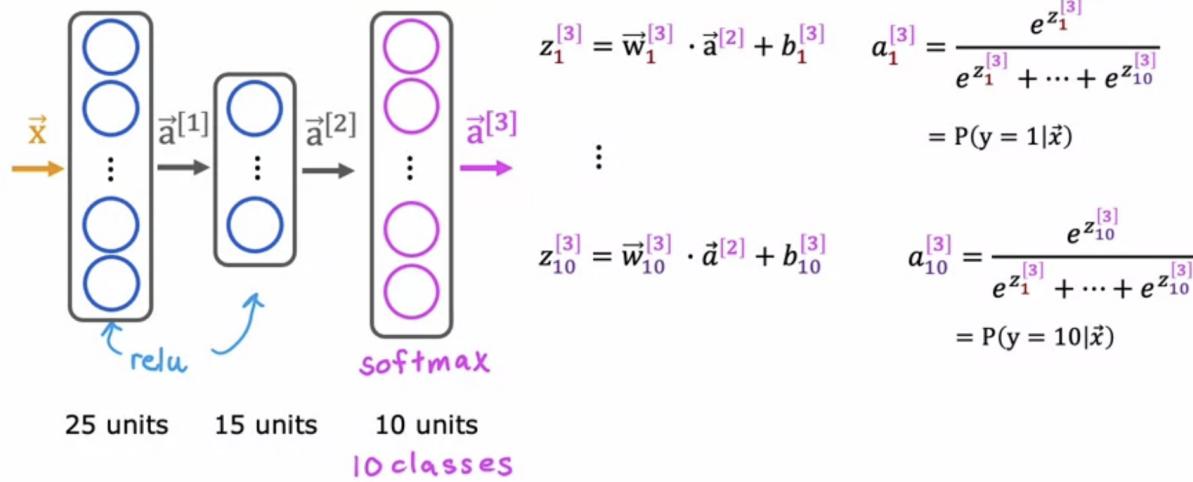
Softmax regression

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

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

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

Neural Network with Softmax output



MNIST with softmax

① specify the model

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

```
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')
])
from tensorflow.keras.losses import
    SparseCategoricalCrossentropy
model.compile(loss= SparseCategoricalCrossentropy() )
model.fit(X, Y, epochs=100)
```

② specify loss and cost

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

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

Numerical Roundoff Errors

More numerically accurate implementation of logistic loss:

$$| + \frac{1}{10,000} \quad | - \frac{1}{10,000}$$

Logistic regression:

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

Original loss

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

model.compile(loss=BinaryCrossEntropy()) ↴

model.compile(loss=BinaryCrossEntropy(from_logits=True)) ↴

More accurate loss (in code)

$$\text{loss} = -y \log\left(\frac{1}{1 + e^{-z}}\right) - (1-y) \log\left(1 - \frac{1}{1 + e^{-z}}\right)$$

More numerically accurate implementation of softmax

Softmax regression

$$(a_1, \dots, a_{10}) = g(z_1, \dots, z_{10})$$

$$\text{Loss} = L(\vec{a}, y) = \begin{cases} -\log a_1 & \text{if } y = 1 \\ \vdots \\ -\log a_{10} & \text{if } y = 10 \end{cases}$$

```
model = Sequential([
    Dense(units=25, activation='relu'),
    Dense(units=15, activation='relu'),
    Dense(units=10, activation='softmax')
])
'linear'  

model.compile(loss=SparseCategoricalCrossEntropy())

```

More Accurate

$$L(\vec{a}, y) = \begin{cases} -\log \frac{e^{z_1}}{e^{z_1} + \dots + e^{z_{10}}} & \text{if } y = 1 \\ \vdots \\ -\log \frac{e^{z_{10}}}{e^{z_1} + \dots + e^{z_{10}}} & \text{if } y = 10 \end{cases}$$

```
model.compile(loss=SparseCategoricalCrossEntropy(from_logits=True))
```

MNIST (more numerically accurate)

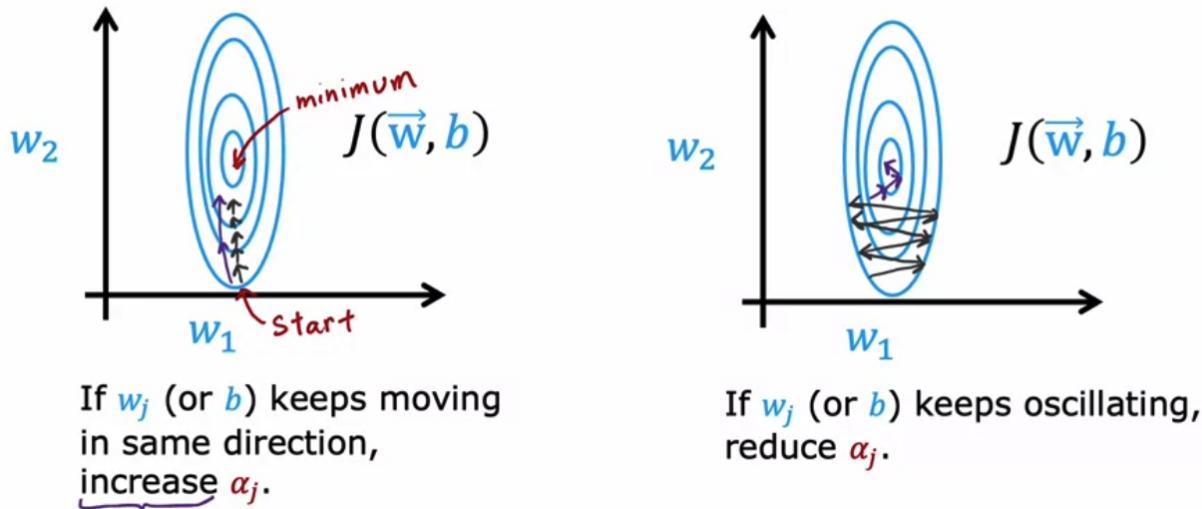
```
model import tensorflow as tf
from tensorflow.keras import Sequential
from tensorflow.keras.layers import Dense
model = Sequential([
    Dense(units=25, activation='relu'),
    Dense(units=15, activation='relu'),
    Dense(units=10, activation='linear') ])
loss from tensorflow.keras.losses import
      SparseCategoricalCrossentropy
model.compile(..., loss=SparseCategoricalCrossentropy(from_logits=True))
fit model.fit(X, Y, epochs=100)
predict logits = model(X)
f_x = tf.nn.softmax(logits)
```

Multi-label classification: single input maps to potentially more than one label. One can build multiple NNs to detect each label or have an output layer with multiple nodes, one for each label, using a sigmoid activation function.

Adam algorithm is an advance of gradient descent, increasing the learning rate if the steps are going in the same direction or decreasing the learning rate if the

steps are going all over the place. It uses different learning rates for each parameter of the model. It's become the standard method.

Adam Algorithm Intuition



MNIST Adam

model

```
model = Sequential([
    tf.keras.layers.Dense(units=25, activation='sigmoid'),
    tf.keras.layers.Dense(units=15, activation='sigmoid'),
    tf.keras.layers.Dense(units=10, activation='linear')
])
```

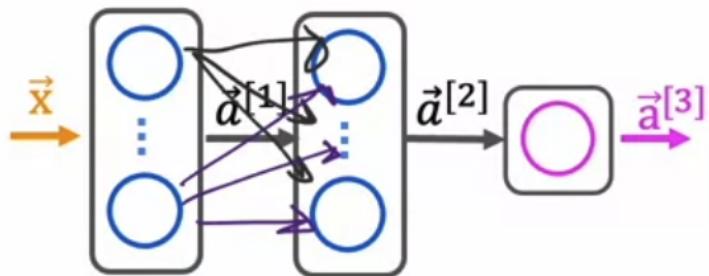
compile

```
model.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=1e-3),
              loss=tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True))
```

fit

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

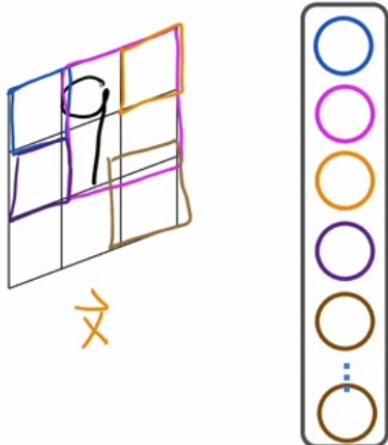
Dense Layer



Each neuron output is a function of
all the activation outputs of the previous layer.

$$\vec{a}_1^{[2]} = g \left(\vec{w}_1^{[2]} \cdot \vec{a}^{[1]} + b_1^{[2]} \right)$$

Convolutional Layer



Each neuron only looks at part of the previous layer's outputs.

Why?

- Faster computation
- Need less training data (less prone to overfitting)

Convolutional Neural Network

