

Neural Networks Basics

Logistic Regression as a Neural Network

$$\hat{y} = \sigma(w^T x + b), \text{ where } \sigma(z) = \frac{1}{1 + e^{-z}}$$

Given $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$, want $\hat{y}^{(i)} \approx y^{(i)}$

Loss(error) function:

Don't use this, non-convex

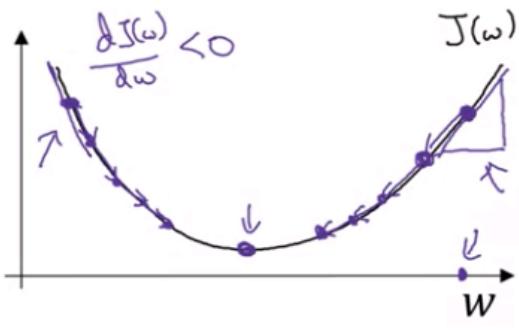
$$\ell(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$$

The loss function computes the error for a single training example; the cost function is the average of the loss functions of the entire training set.

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$

Want to find w, b that minimize $J(w, b)$

Gradient Descent



Repeat {
 $w := w - \alpha$
} learning rate
 $\boxed{\frac{dJ(w)}{dw}}$
 $w := w - \alpha \underline{\frac{dJ(w)}{dw}}$
 $\underbrace{\frac{dJ(w)}{dw}}_? = ?$

$$J(w, b)$$

$$w := w - \alpha \boxed{\frac{dJ(w, b)}{dw}}$$

$$\boxed{\frac{\partial J(w, b)}{\partial w}}$$

2

$$b := b - \alpha \frac{dJ(w, b)}{db}$$

Computation Graph

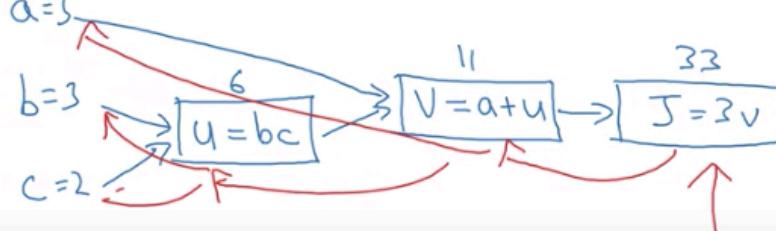
$$J(a, b, c) = 3(a + bc) = 3(5 + 3 \cdot 2) = 33$$

$$\underbrace{\begin{array}{c} u \\ v \end{array}}_{J}$$

$$u = bc$$

$$v = a + u$$

$$J = 3v$$



Computing derivatives

$\frac{\partial J}{\partial a} \rightarrow \underline{d_a = 3}$

$\frac{\partial J}{\partial b} \rightarrow \underline{d_b = 6}$

$\frac{\partial J}{\partial c} \rightarrow \underline{d_c = 9}$

$\frac{\partial J}{\partial u} = 3 = \frac{\partial J}{\partial v} \cdot \frac{\partial v}{\partial u}$

$\frac{\partial J}{\partial b} = \frac{\partial J}{\partial v} \cdot \frac{\partial v}{\partial u} \cdot \frac{\partial u}{\partial b} = 6$

$\frac{\partial J}{\partial b} = \frac{\partial J}{\partial v} \cdot \frac{\partial v}{\partial u} \cdot \frac{\partial u}{\partial c} \cdot \frac{\partial c}{\partial b} = 6$

$v = a + u \quad \underline{d_v = 3} \quad \underline{\frac{\partial J}{\partial v}}$

$J = 3v \quad \uparrow$

$u = bc \quad \underline{d_u = 3}$

$a = 5 \rightarrow \underline{d_a = 3}$

$b = 3 \rightarrow \underline{d_b = 6}$

$c = 2 \rightarrow \underline{d_c = 9}$

$u = 6 \rightarrow \underline{6.001}$

$v = 11 \rightarrow \underline{11.001}$

$J = 33 \rightarrow \underline{33.003}$

$b = 3 \rightarrow \underline{3.001}$

$u = b \cdot c = 6 \rightarrow \underline{6.002}$

$J = 33.006 \quad c = 2 \quad .006$

$v = 11.002 \quad J = 33V$

13:32 / 14:33

Gradient Descent

Logistic regression derivatives

x_1

w_1

x_2

w_2

b

$z = w_1x_1 + w_2x_2 + b \rightarrow a = \sigma(z) \rightarrow \mathcal{L}(a, y)$

$\frac{\partial z}{\partial z} = \frac{\partial \mathcal{L}}{\partial z} = \frac{\partial \mathcal{L}(ay)}{\partial z}$

$\frac{\partial a}{\partial z} = \frac{\partial \mathcal{L}(ay)}{\partial a}$

$= \frac{\partial \mathcal{L}}{\partial a} \cdot \frac{\partial a}{\partial z}$

$= -\frac{y}{a} + \frac{1-y}{1-a}$

$\frac{\partial \mathcal{L}}{\partial w_1} = "dw_1" = x_1 \cdot dz \quad dz = x_2 \cdot dz \quad db = dz$

Logistic regression on m examples

$$\begin{aligned} \mathcal{J}(\omega, b) &= \frac{1}{m} \sum_{i=1}^m \ell(a^{(i)}, y^{(i)}) \\ \Rightarrow a^{(i)} &= \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(\omega^\top x^{(i)} + b) \end{aligned} \quad (\underline{x^{(i)}}, \underline{y^{(i)}})$$

$$\underline{d\omega_1}, \underline{d\omega_2}, \underline{db}$$

$$\underline{\frac{\partial}{\partial \omega_1} \mathcal{J}(\omega, b)} = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial}{\partial \omega_1} \ell(a^{(i)}, y^{(i)})}_{\underline{d\omega_1}} - (\underline{x^{(i)}}, \underline{y^{(i)}})$$

Step 1 : $\frac{dL}{da}$

$$L = -(y \times \log(a) + (1 - y) \times \log(1 - a))$$

$$\frac{dL}{da} = -y \times \frac{1}{a} - (1 - y) \times \frac{1}{1-a} \times -1$$

$$\begin{aligned} \frac{d}{dx} \sigma(x) &= \frac{d}{dx} \left[\frac{1}{1 + e^{-x}} \right] \\ &= \frac{d}{dx} (1 + e^{-x})^{-1} \\ &= -(1 + e^{-x})^{-2} (-e^{-x}) \\ &= \frac{e^{-x}}{(1 + e^{-x})^2} \\ &= \frac{1}{1 + e^{-x}} \cdot \frac{e^{-x}}{1 + e^{-x}} \\ &= \frac{1}{1 + e^{-x}} \cdot \frac{(1 + e^{-x}) - 1}{1 + e^{-x}} \\ &= \frac{1}{1 + e^{-x}} \cdot \left(\frac{1 + e^{-x}}{1 + e^{-x}} - \frac{1}{1 + e^{-x}} \right) \\ &= \frac{1}{1 + e^{-x}} \cdot \left(1 - \frac{1}{1 + e^{-x}} \right) \\ &= \sigma(x) \cdot (1 - \sigma(x)) \end{aligned}$$

In the previous video, Andrew refers to $dz = a(1 - a)$

Note that Andrew is using "dz" as a shorthand to refer to $\frac{da}{dz} = a(1 - a)$.

To clarify, earlier in this week's videos, Andrew used the name "dz" to refer to a different derivative: $\frac{dL}{dz} = a - y$.

Recall that the relationship between $\frac{dL}{dz}$ and $\frac{da}{dz}$ is:

$$\frac{dL}{dz} = \frac{dL}{da} \times \frac{da}{dz}$$

$$\frac{dL}{dz} = \frac{a-y}{a(1-a)} \times a(1-a) = a - y$$

Vectorization

Logistic regression derivatives

$$\begin{aligned}
 J &= 0, \quad \boxed{\cancel{dw_1 = 0}, \cancel{dw_2 = 0}}, \quad db = 0 \quad dw = np.zeros((n_x, 1)) \\
 \rightarrow \text{for } i &= 1 \text{ to } m: \\
 z^{(i)} &= w^T x^{(i)} + b \\
 a^{(i)} &= \sigma(z^{(i)}) \\
 J &+= -[y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})] \\
 \downarrow \text{for } j &= 1 \dots n_x \\
 dz^{(i)} &= a^{(i)}(1 - a^{(i)}) \\
 \boxed{\begin{array}{l|l} dw_1 += x_1^{(i)} dz^{(i)} & n_x = 2 \\ dw_2 += x_2^{(i)} dz^{(i)} & \\ \end{array}} & dw += x^{(i)} dz^{(i)} \\
 db &+= dz^{(i)} \\
 J &= J/m, \quad \boxed{dw_1 = dw_1/m, \quad dw_2 = dw_2/m, \quad db = db/m} \quad dw / = m
 \end{aligned}$$

Vectorizing Logistic Regression

$$\begin{aligned}
 dz^{(1)} &= a^{(1)} - y^{(1)} & dz^{(2)} &= a^{(2)} - y^{(2)} & \dots \\
 dz &= \boxed{[dz^{(1)} \ dz^{(2)} \ \dots \ dz^{(m)}]}_{{1 \times m}} \leftarrow & & & \\
 A &= [a^{(1)} \ \dots \ a^{(m)}], \quad Y = [y^{(1)} \ \dots \ y^{(m)}] & & & \\
 \rightarrow dz &= A - Y = [a^{(1)} - y^{(1)} \ a^{(2)} - y^{(2)} \ \dots] & & & \\
 \rightarrow \boxed{\begin{array}{l} dw = 0 \\ dw += X^{(1)} dz^{(1)} \\ dw += X^{(2)} dz^{(2)} \\ \vdots \\ dw / = m \end{array}} & \left| \begin{array}{l} db = 0 \\ db += dz^{(1)} \\ db += dz^{(2)} \\ \vdots \\ db / = m \end{array} \right. & & \begin{aligned} db &= \frac{1}{m} \sum_{i=1}^m dz^{(i)} \\ &= \frac{1}{m} \underline{\text{np.sum}(dz)} \\ dw &= \frac{1}{m} X dz^T \\ &= \frac{1}{m} \left[\begin{array}{c|c} \underline{x^{(1)}} & \dots & \underline{x^{(m)}} \\ \hline 1 & \dots & 1 \end{array} \right] \left[\begin{array}{c} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{array} \right] \\ &= \frac{1}{m} \left[\underline{x^{(1)} dz^{(1)}} + \dots + \underline{x^{(m)} dz^{(m)}} \right]_{n \times 1} \end{aligned} & &
 \end{aligned}$$

Implementing Logistic Regression

```

 $J = 0, dw_1 = 0, dw_2 = 0, db = 0$ 
for i = 1 to m:
     $z^{(i)} = w^T x^{(i)} + b \leftarrow$ 
     $a^{(i)} = \sigma(z^{(i)}) \leftarrow$ 
     $J += -[y^{(i)} \log a^{(i)} + (1 - y^{(i)}) \log(1 - a^{(i)})]$ 
     $dz^{(i)} = a^{(i)} - y^{(i)} \leftarrow$ 
     $\begin{cases} dw_1 += x_1^{(i)} dz^{(i)} \\ dw_2 += x_2^{(i)} dz^{(i)} \\ db += dz^{(i)} \end{cases}$ 
 $J = J/m, dw_1 = dw_1/m, dw_2 = dw_2/m$ 
 $db = db/m$ 

```

$$\begin{aligned}
z &= w^T X + b \\
&= n \cdot \text{dot}(w.T, X) + b \\
A &= \sigma(z) \\
d\bar{z} &= A - Y \\
dw &= \frac{1}{m} X d\bar{z}^T \\
db &= \frac{1}{m} \text{np.sum}(d\bar{z}) \\
w &:= w - \alpha dw \\
b &:= b - \alpha db
\end{aligned}$$

Andrew

Logistic regression cost function

$$\begin{aligned}
&\rightarrow \boxed{\text{If } y = 1: p(y|x) = \hat{y}} \\
&\rightarrow \boxed{\text{If } y = 0: p(y|x) = 1 - \hat{y}}
\end{aligned}
\quad \left. \begin{array}{l} \{ \\ \} \end{array} \right\} p(y|x)$$

$$p(y|x) = \underbrace{\hat{y}^y}_{\text{If } y=1} \underbrace{(1-\hat{y})^{(1-y)}}_{\text{If } y=0}$$

$$\begin{aligned}
\text{If } y=1: p(y|x) &= \hat{y}^y \underbrace{(1-\hat{y})^0}_{=1} \\
\text{If } y=0: p(y|x) &= \hat{y}^0 \underbrace{(1-\hat{y})^{(1-y)}}_{=1} = 1 \times (1-\hat{y}) = 1 - \hat{y}
\end{aligned}$$

$$\begin{aligned}
\uparrow \log p(y|x) &= \log \hat{y}^y (1-\hat{y})^{(1-y)} = y \log \hat{y} + (1-y) \log(1-\hat{y}) \\
&= \overbrace{-\ell(\hat{y}, y)}^{\text{J}}
\end{aligned}$$

Andrew Ng

Cost on m examples

$$\begin{aligned}
 \log p(\text{labels in training set}) &= \log \prod_{i=1}^m p(y^{(i)} | x^{(i)}) \leftarrow \\
 \log p(\text{---}) &= \sum_{i=1}^m \underbrace{\log p(y^{(i)} | x^{(i)})}_{-L(\hat{y}^{(i)}, y^{(i)})} \\
 &= -\sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) \\
 \text{Cost: } \underset{\text{(minimize)}}{\overline{J(w, b)}} &= \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})
 \end{aligned}$$

Maximum likelihood estimation ↗

Code

For convenience, you should now reshape images of shape (num_px, num_px, 3) in a numpy-array of shape (num_px * num_px * 3, 1). After this, our training (and test) dataset is a numpy-array where each column represents a flattened image. There should be m_train (respectively m_test) columns.

Exercise: Reshape the training and test data sets so that images of size (num_px, num_px, 3) are flattened into single vectors of shape (num_px * num_px * 3, 1).

A trick when you want to flatten a matrix X of shape (a,b,c,d) to a matrix X_flatten of shape (b*c*d, a) is to use:

```
x_flatten = x.reshape(x.shape[0], -1).T      # x.T is the transpose of x
```

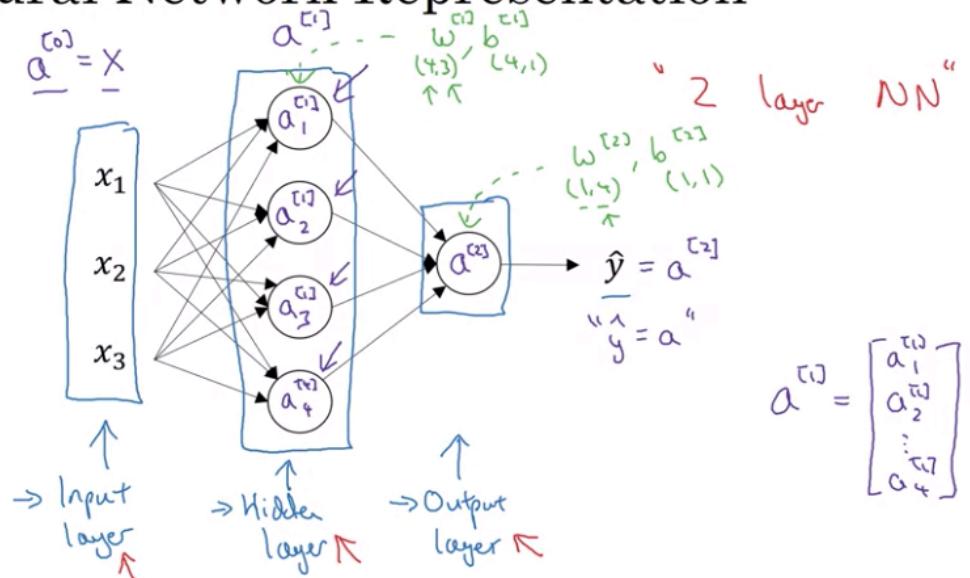
To represent color images, the red, green and blue channels (RGB) must be specified for each pixel, and so the pixel value is actually a vector of three numbers ranging from 0 to 255.

One common preprocessing step in machine learning is to center and standardize your dataset, meaning that you subtract the mean of the whole numpy array from each example, and then divide each example by the standard deviation of the whole numpy array. But for picture datasets, it is simpler and more convenient and works almost as well to just divide every row of the dataset by 255 (the maximum value of a pixel channel).

Let's standardize our dataset.

Neural Network

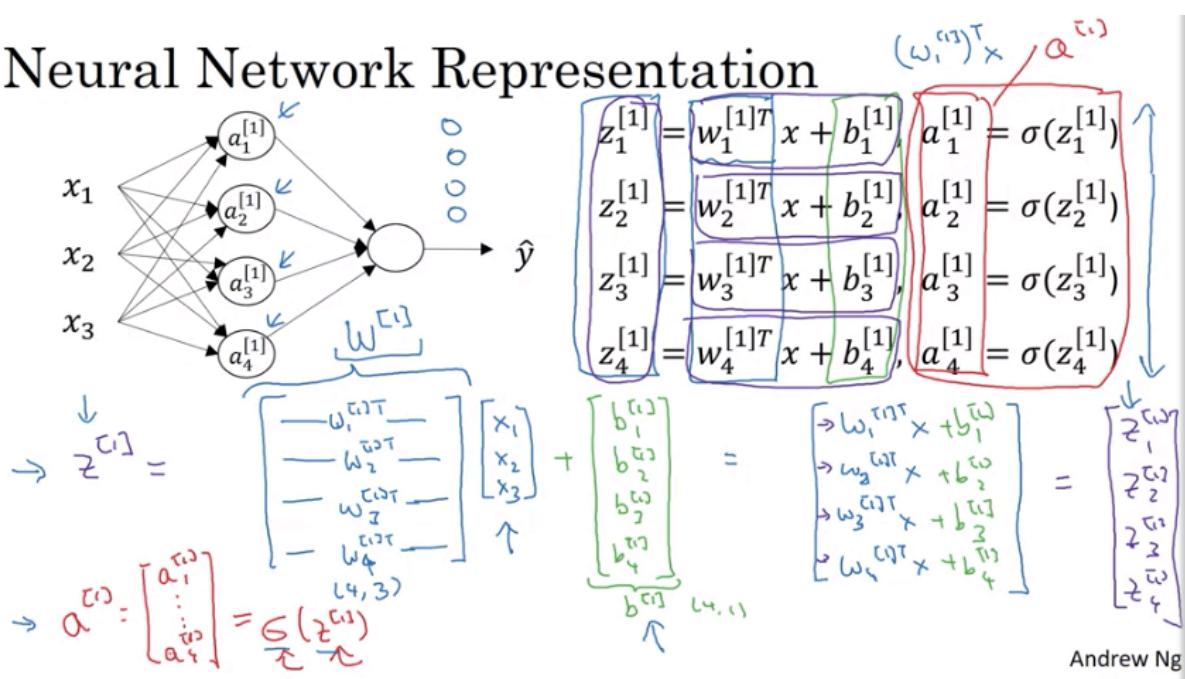
Neural Network Representation



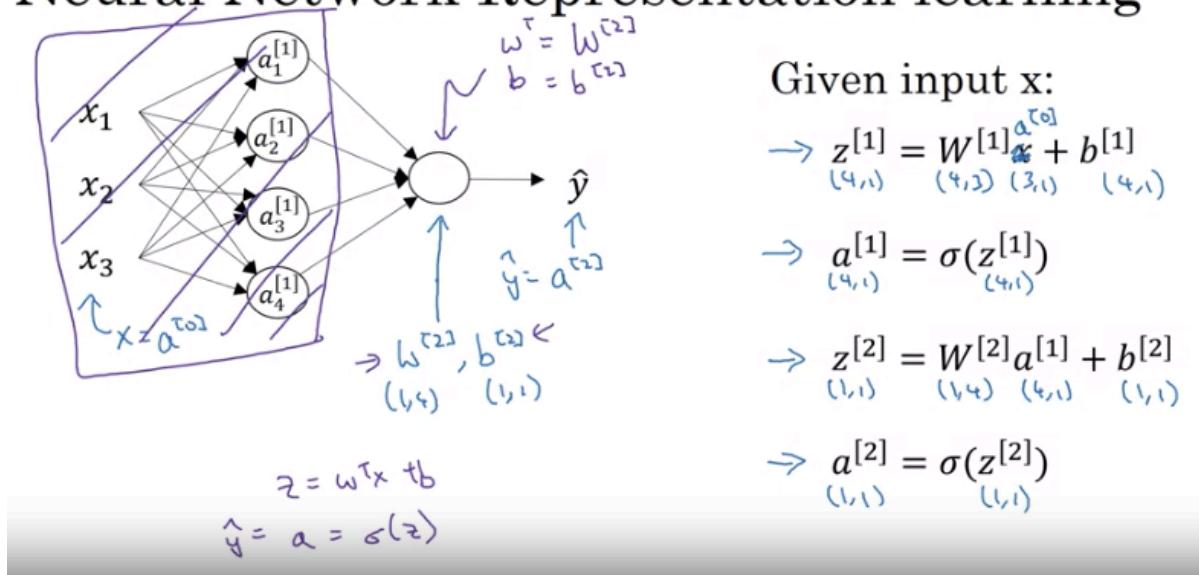
Shallow neural networks

Representation

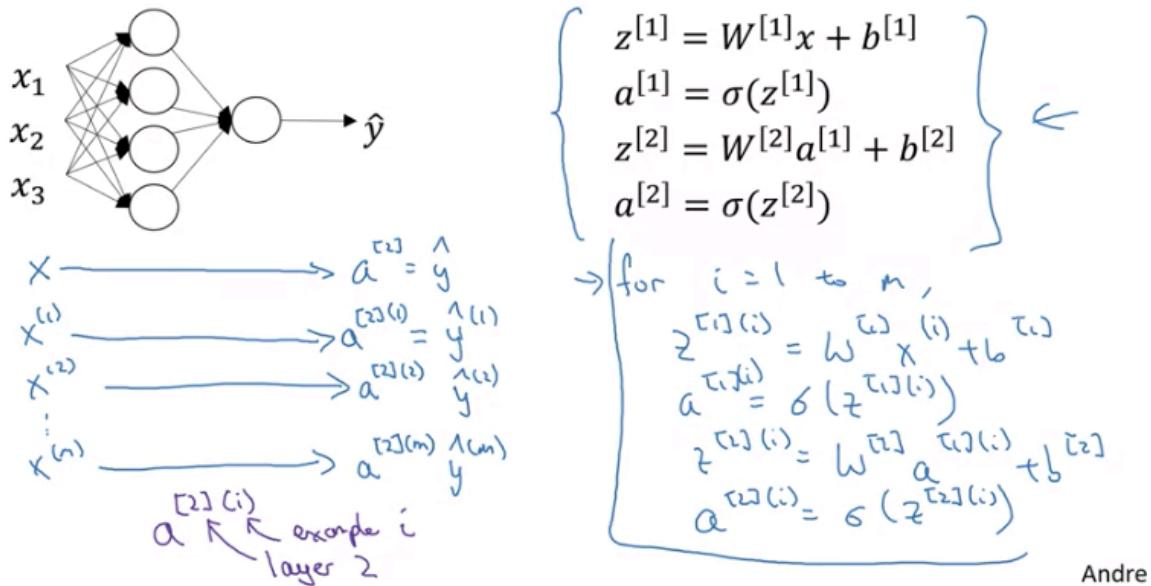
Neural Network Representation



Neural Network Representation learning



Vectorizing across multiple examples



Vectorizing across multiple examples

for $i = 1$ to m :

$$z^{[1](i)} = W^{[1]}x^{(i)} + b^{[1]}$$

$$a^{[1](i)} = \sigma(z^{[1](i)})$$

$$z^{[2](i)} = W^{[2]}a^{[1](i)} + b^{[2]}$$

$$a^{[2](i)} = \sigma(z^{[2](i)})$$

$$X = \begin{bmatrix} & & & \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ & & & \end{bmatrix}$$

\uparrow \uparrow
 (n_x, m)

↑ ↑
hidden units.

↑ ↑
training examples

$$\begin{aligned} z^{[1]} &= W^{[1]}X + b^{[1]} \\ \Rightarrow A^{[1]} &= \sigma(z^{[1]}) \\ \Rightarrow z^{[2]} &= W^{[2]}A^{[1]} + b^{[2]} \\ \Rightarrow A^{[2]} &= \sigma(z^{[2]}) \end{aligned}$$

$$\begin{aligned} z^{[1]} &= \begin{bmatrix} z^{1} & z^{[1](2)} & \dots & z^{[1](m)} \end{bmatrix} \\ A^{[1]} &= \begin{bmatrix} a^{1} & a^{[1](2)} & \dots & a^{[1](m)} \end{bmatrix} \\ z^{[2]} &= \begin{bmatrix} z^{[2](1)} & z^{2} & \dots & z^{[2](m)} \end{bmatrix} \\ A^{[2]} &= \begin{bmatrix} a^{[2](1)} & a^{2} & \dots & a^{[2](m)} \end{bmatrix} \end{aligned}$$

Andrew Ng

X: 竖方向 features

横方向: training example

Justification for vectorized implementation

$$\begin{aligned}
 z^{(1)(1)} &= w^{(1)} x^{(1)} + b^{(1)}, & z^{(1)(2)} &= w^{(1)} x^{(1)} + b^{(1)}, & z^{(1)(3)} &= w^{(1)} x^{(1)} + b^{(1)} \\
 w^{(1)} &= \begin{bmatrix} \cdot & \cdot & \cdot \end{bmatrix}, & w^{(1)} x^{(1)} &= \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}, & w^{(1)} x^{(1)} &= \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}, & w^{(1)} x^{(1)} &= \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \\
 w^{(1)} \begin{bmatrix} 1 & x^{(1)} & x^{(2)} & \dots \end{bmatrix} &= \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix} = \begin{bmatrix} z^{(1)(1)} \\ z^{(1)(2)} \\ z^{(1)(3)} \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} = \begin{bmatrix} z^{(1)} \\ z^{(1)} \\ z^{(1)} \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} = z^{(1)} & w^{(1)} x^{(1)} &= z^{(1)(1)}, & w^{(1)} x^{(1)} &= z^{(1)(2)}, & w^{(1)} x^{(1)} &= z^{(1)(3)}
 \end{aligned}$$

Recap of vectorizing across multiple examples

Diagram of a neural network with input x_1, x_2, x_3 and output \hat{y} .

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix}$$

$$A^{[1]} = \begin{bmatrix} a^{1} & a^{[1](2)} & \dots & a^{[1](m)} \end{bmatrix}$$

Code snippets:

```

for i = 1 to m
    z^{[1](i)} = W^{[1]}x^{(i)} + b^{[1]}
    → a^{[1](i)} = σ(z^{[1](i)})
    → z^{[2](i)} = W^{[2]}a^{[1](i)} + b^{[2]}
    → a^{[2](i)} = σ(z^{[2](i)})
```

```

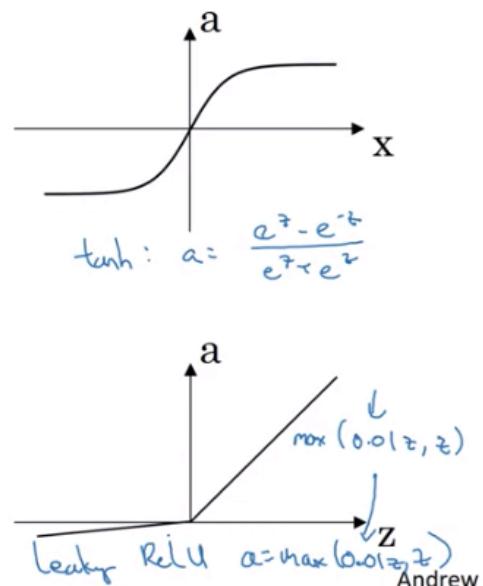
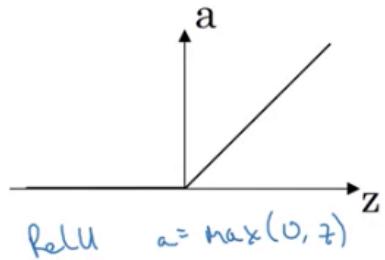
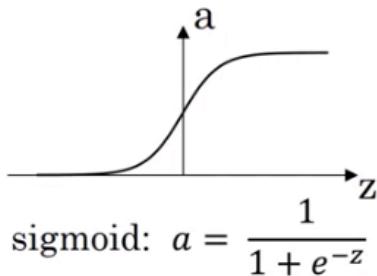
Z^{[1]} = W^{[1]}X + b^{[1]} ← w^{[1]}A^{[1]} + b^{[1]}
A^{[1]} = σ(Z^{[1]})
```

```

Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]} ← w^{[2]}A^{[1]} + b^{[2]}
A^{[2]} = σ(Z^{[2]})
```

Andrew Ng

Pros and cons of activation functions

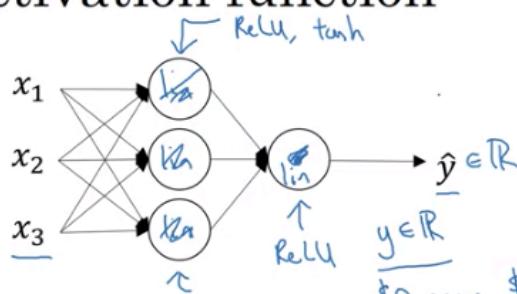


Activation function:

1. Sigmoid: 如果 output 是 binary classification, 可以考慮 sigmoid function, 一般 never used
2. tanh function:
3. ReLU Rectified Linear Unit $a = \max(0, z)$: by default, 很多人用, faster, slope 1, or 0
4. Leaky ReLU

Why need non-linear activation functions

Activation function



Given x :

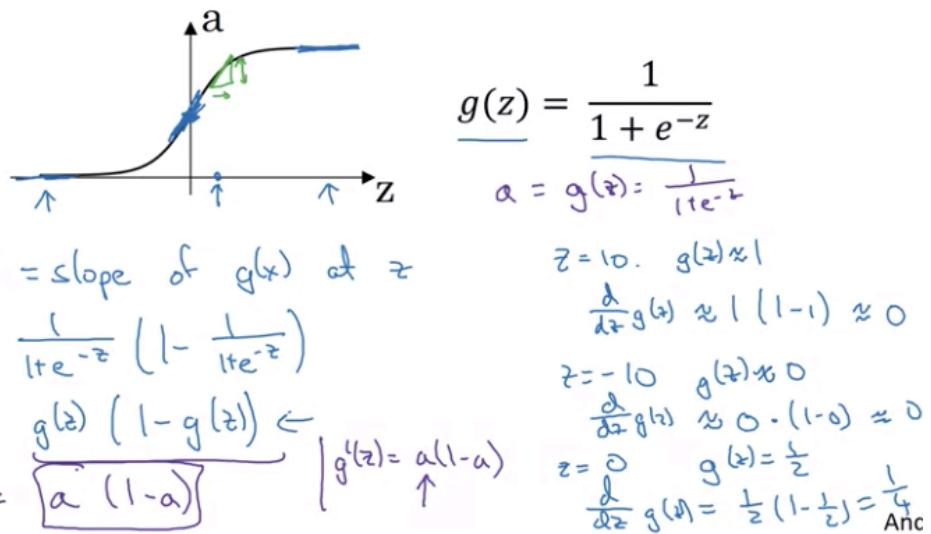
$$\begin{aligned} \rightarrow z^{[1]} &= W^{[1]}x + b^{[1]} \\ \rightarrow a^{[1]} &= g^{[1]}(z^{[1]}) = z^{[1]} \quad "linear\ activation\ function" \\ \rightarrow z^{[2]} &= W^{[2]}a^{[1]} + b^{[2]} \\ \rightarrow a^{[2]} &= g^{[2]}(z^{[2]}) = z^{[2]} \end{aligned}$$

$$\begin{aligned} a^{[1]} &= z^{[1]} = \underbrace{W^{[1]}x + b^{[1]}}_{a^{[1]}} \\ a^{[2]} &= z^{[2]} = \underbrace{W^{[2]}a^{[1]} + b^{[2]}}_{a^{[2]}} \\ a^{[1]} &= \underbrace{W^{[1]}(W^{[1]}x + b^{[1]})}_{a^{[1]}} + b^{[2]} \\ &= (\underbrace{W^{[2]}W^{[1]}}_{w'})x + (\underbrace{W^{[2]}b^{[1]}}_{b'}) \\ &= w'x + b' \\ g(w') &= w' \end{aligned}$$

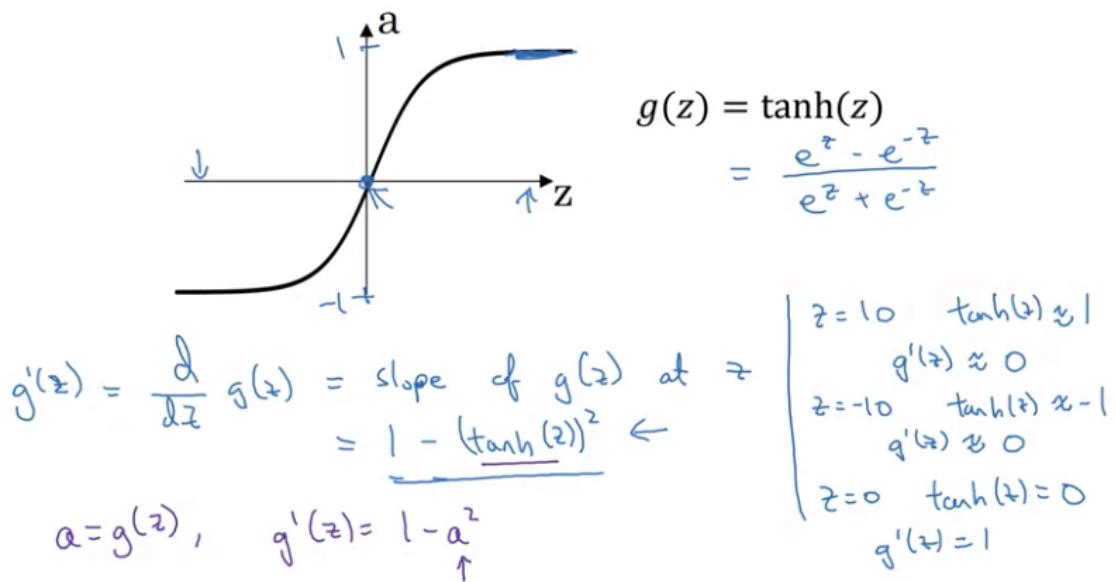
Andrew Ng

Derivatives of activation functions

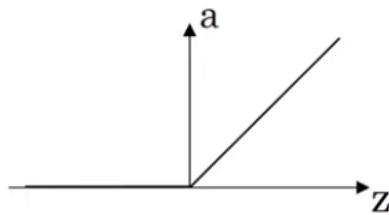
Sigmoid activation function



Tanh activation function



ReLU and Leaky ReLU



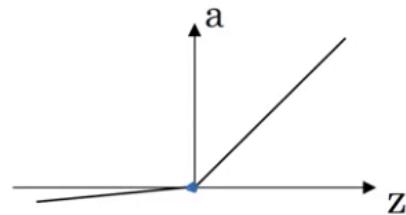
ReLU

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

$$\rightarrow g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

~~unlike if $z=0$~~

$z = 0.0000000000$



Leaky ReLU

$$g(z) = \max(0.01z, z)$$

$$g'(z) = \begin{cases} 0.01 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

Andrei

Gradient descent for Neural Networks

`np.sum(keepdims=True)`: prevent rank one arrays

Gradient descent for neural networks

Parameters: $\underline{w^{[0]}}$, $\underline{b^{[0]}}$, $\underline{w^{[1]}}$, $\underline{b^{[1]}}$, $\underline{w^{[2]}}$, $\underline{b^{[2]}}$

$$n_x = n^{[0]}, n^{[1]}, \underline{n^{[2]} = 1}$$

Cost function: $J(\underline{w^{[0]}}, \underline{b^{[0]}}, \underline{w^{[1]}}, \underline{b^{[1]}}, \underline{w^{[2]}}, \underline{b^{[2]}}) = \frac{1}{m} \sum_{i=1}^m \sum_{a^{[2]}} L(\hat{y}, y)$

Gradient Descent:

\rightarrow Repeat {
 → Compute predictions $(\hat{y}^{(i)}, i=1 \dots m)$
 $\underline{\Delta w^{[0]}} = \frac{\partial J}{\partial w^{[0]}}, \underline{\Delta b^{[0]}} = \frac{\partial J}{\partial b^{[0]}}, \dots$
 $w^{[0]} := w^{[0]} - \alpha \underline{\Delta w^{[0]}}$
 $b^{[0]} := b^{[0]} - \alpha \underline{\Delta b^{[0]}}$

Formulas for computing derivatives

Forward propagation:

$$z^{(1)} = w^{(1)} X + b^{(1)}$$

$$A^{(1)} = g^{(1)}(z^{(1)}) \leftarrow$$

$$z^{(2)} = w^{(2)} A^{(1)} + b^{(2)}$$

$$A^{(2)} = g^{(2)}(z^{(2)}) = \underline{g}(z^{(2)})$$

Back propagation:

$$\delta z^{(2)} = A^{(2)} - Y \leftarrow$$

$$\delta w^{(2)} = \frac{1}{m} \delta z^{(2)} A^{(1)T}$$

$$\delta b^{(2)} = \frac{1}{m} \text{np.sum}(\delta z^{(2)}, \text{axis}=1, \text{keepdims=True})$$

$$\delta z^{(1)} = \underbrace{(w^{(1)T} \delta z^{(2)})}_{(n^{(2)}, m)} \times \underbrace{g'(z^{(1)})}_{\text{element-wise product}} \quad (n^{(1)}, m)$$

$$\delta w^{(1)} = \frac{1}{m} \delta z^{(1)} X^T$$

$$\delta b^{(1)} = \frac{1}{m} \text{np.sum}(\delta z^{(1)}, \text{axis}=1, \text{keepdims=True})$$

$$Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]$$

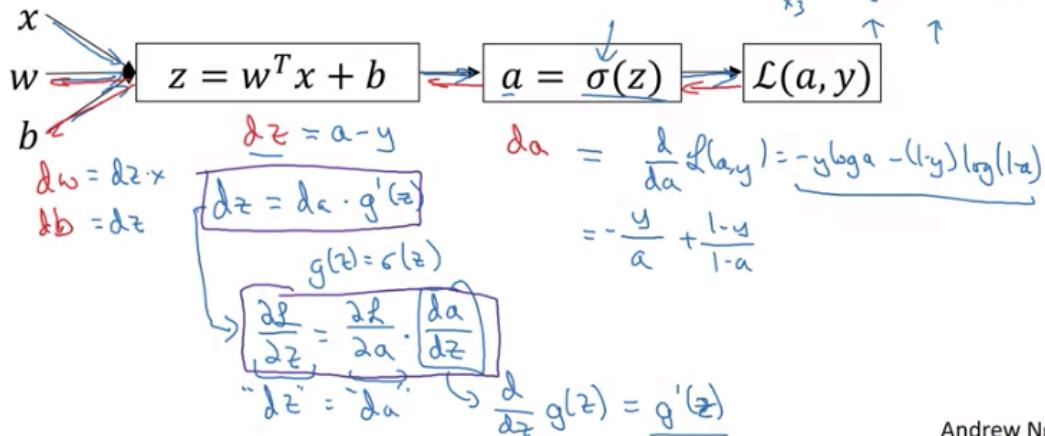
$$(n^{(0)}) \leftarrow$$

$$\downarrow (n^{(1)}, 1) \leftarrow$$

Backpropagation intuition

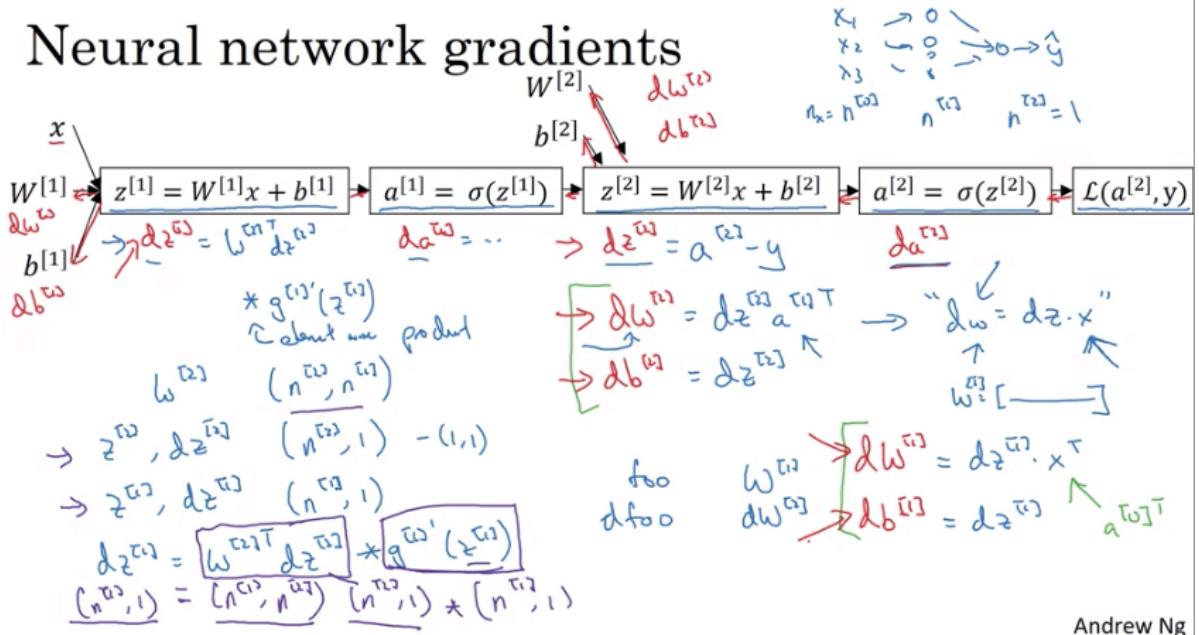
Computing gradients

Logistic regression



Andrew Ng

Neural network gradients



Andrew Ng

<https://medium.com/@pdquant/all-the-backpropagation-derivatives-d5275f727f60>

Summary of gradient descent

$dz^{[2]} = a^{[2]} - y$	$dZ^{[2]} = A^{[2]} - Y$	$J(\cdot) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}_i, y_i)$
$dW^{[2]} = dz^{[2]} a^{[1]T}$	$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$	
$db^{[2]} = dz^{[2]}$	$db^{[2]} = \frac{1}{m} np.sum(dZ^{[2]}, axis=1, keepdims=True)$	
$dz^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]}'(z^{[1]})$	$dZ^{[1]} = \underbrace{W^{[2]T} dZ^{[2]}}_{(n^{[2]}, m)} * \underbrace{g^{[1]}'(Z^{[1]})}_{(n^{[1]}, m)}$	elementwise product
$dW^{[1]} = dz^{[1]} x^T$	$dW^{[1]} = \frac{1}{m} dZ^{[1]} X^T$	
$db^{[1]} = dz^{[1]}$	$db^{[1]} = \frac{1}{m} np.sum(dZ^{[1]}, axis=1, keepdims=True)$	

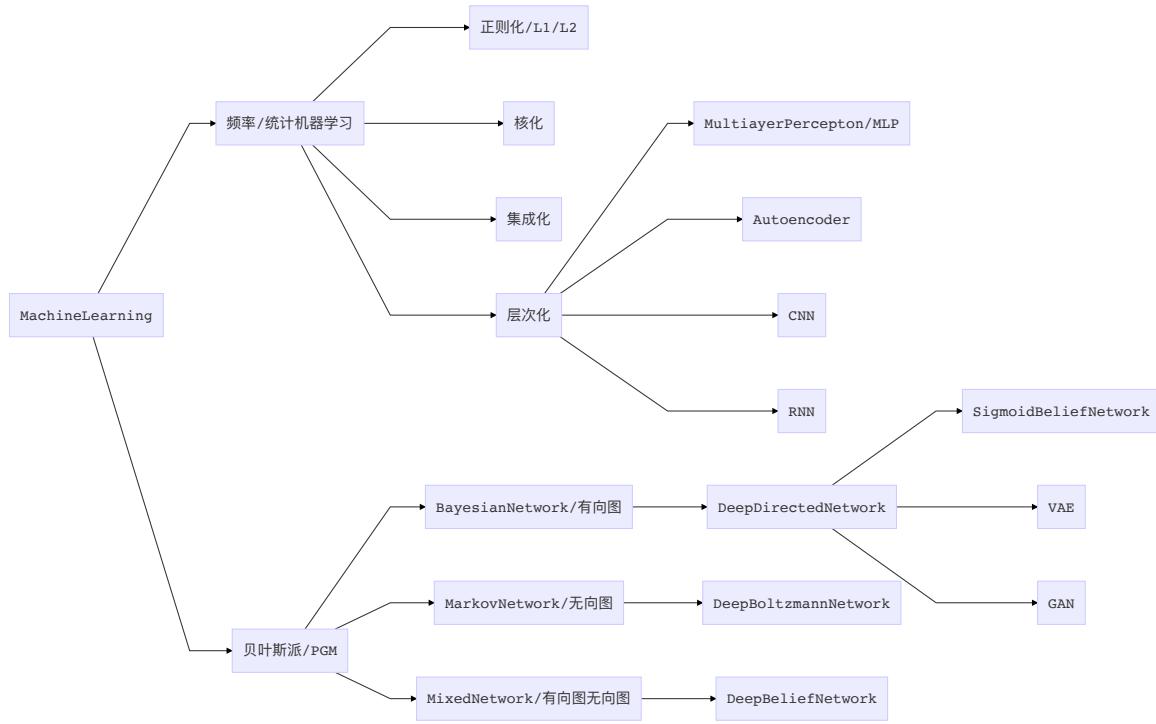
<https://www.coursera.org/learn/neural-networks-deep-learning/discussions/weeks/3/threads/a38VuhyMEei5zw6yFhWyQg>

Random Initialization

LR can be initialized as 0, but not neural network

Supplement

基础



统计机器学习

- 正则化: Loss Function + regularizer(L1/L2)
- 核化: Kernel SVM
- 集成化: AdaBoost, RandomForest
- 层次化: Neural Network/Deep Neural Network
 - MLP(MultiLayer Perceptron)
 - Autoencoder
 - CNN
 - RNN

贝叶斯派

- BayesianNetwork \Rightarrow Deep Directed Network
 - Sigmoid Belief Network
 - Variational Autoencoder(VAE)

- GAN
- Markov Network \Rightarrow Deep Boltzmann Network
- Mixed Netowrk \Rightarrow Deep Belief Network

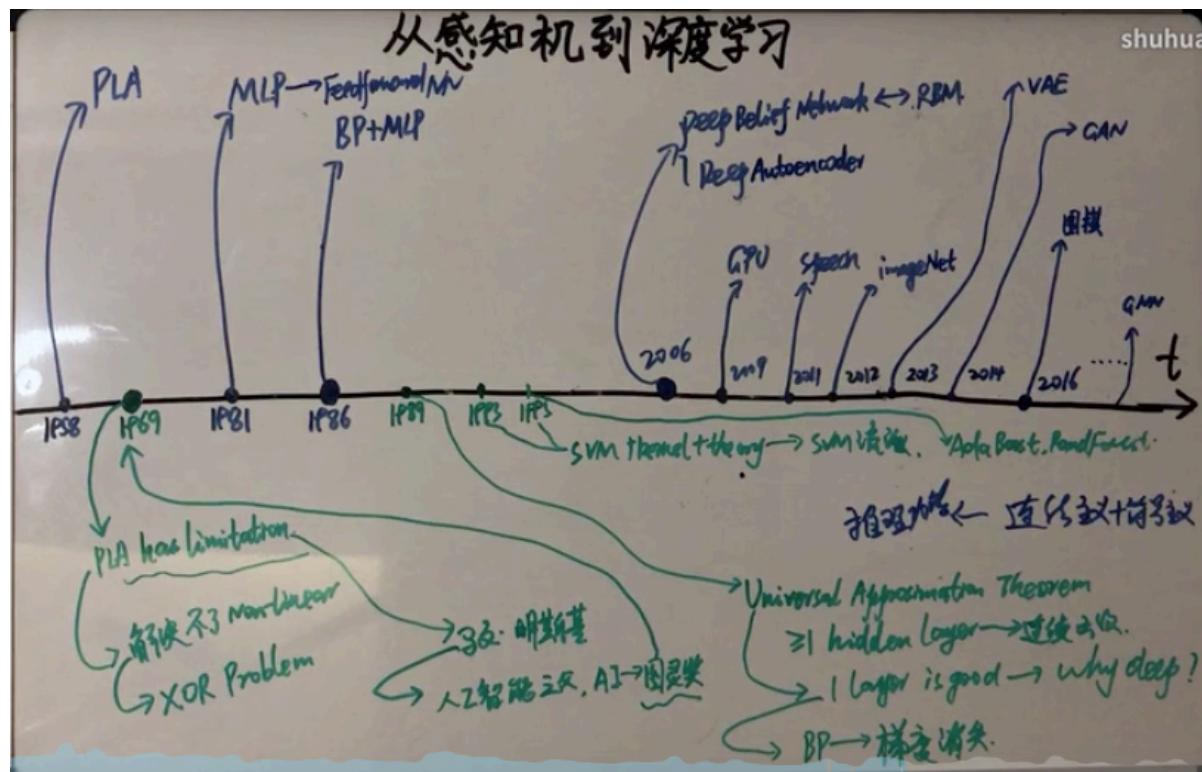
上述都是：Deep Generative Model

狭义的DeepLearning: Deep Neural Network

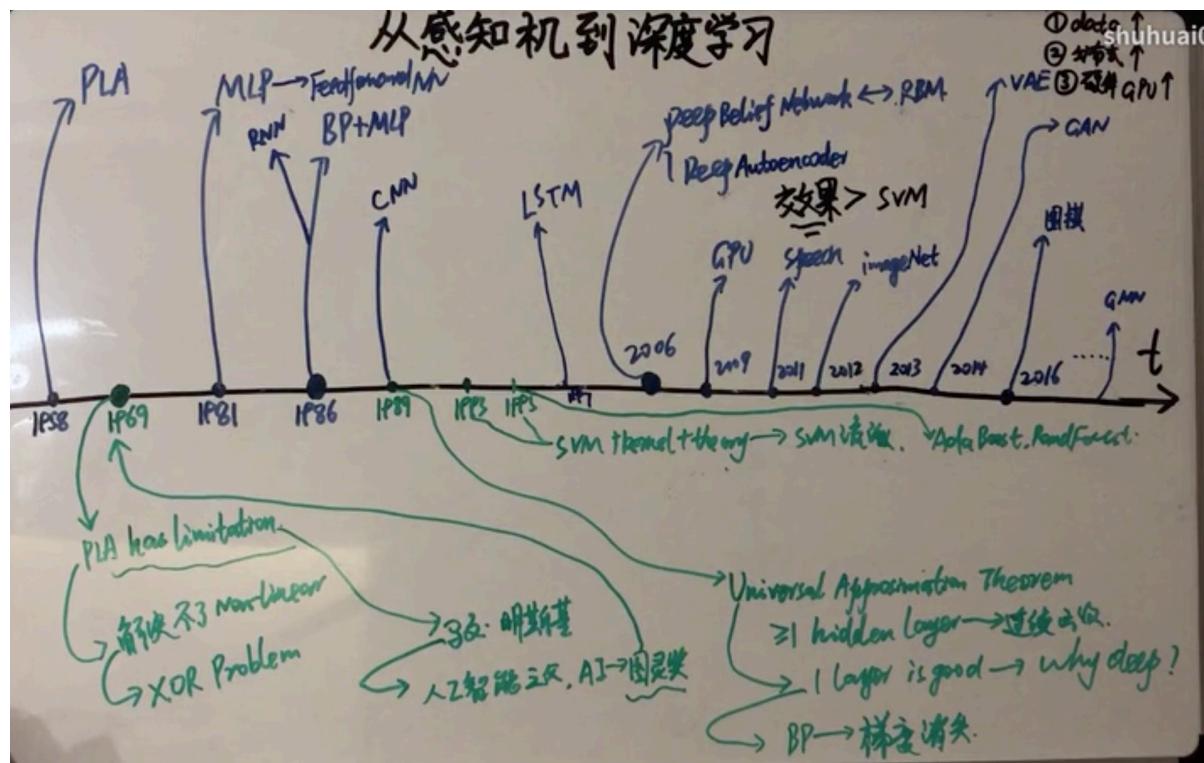
其实应该包括：

1. Deep Neural Network
2. Deep Generative Model (当层次非常多，推断非常困难)

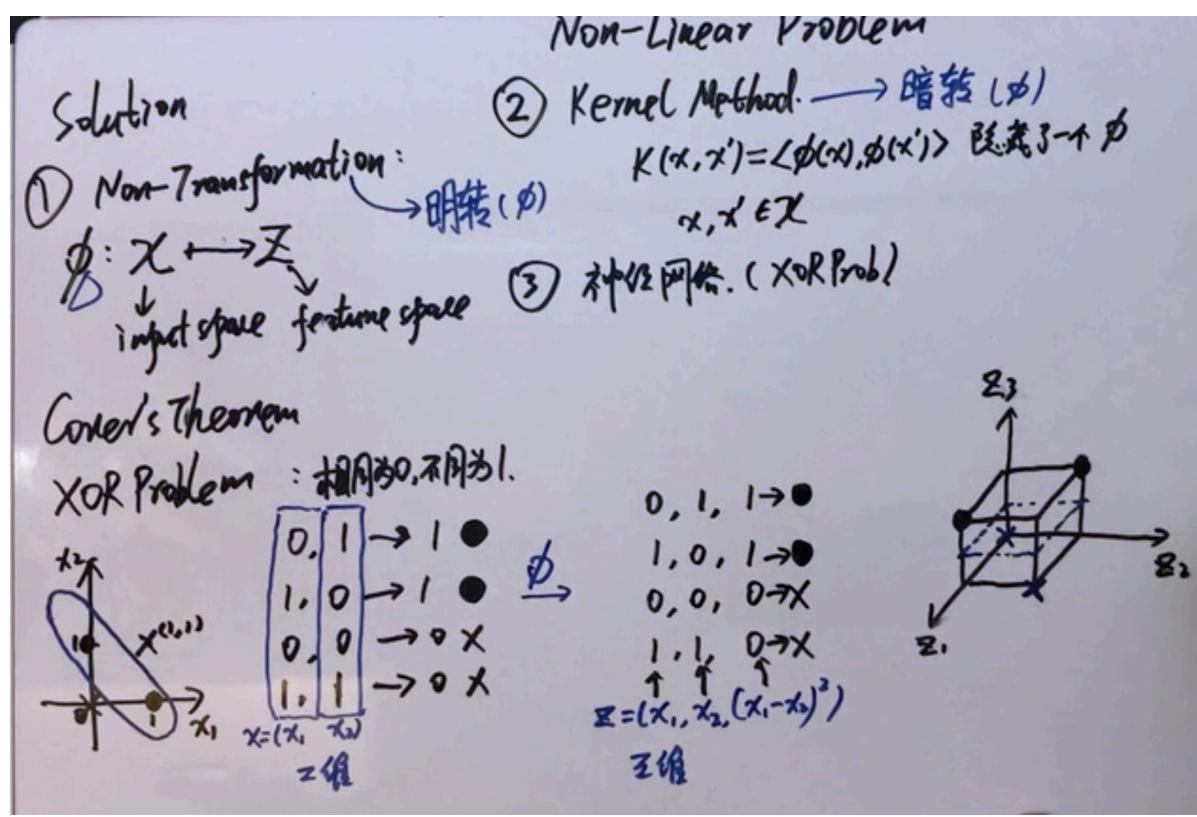
时间线



1. 深度学习的理论在2006年已经成型，直到现在，理论并没有根本性突破。
2. 为什么take off
 1. data
 2. 分布式
 3. 硬件 GPU
 4. 效果



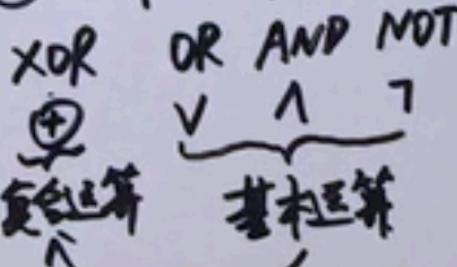
Non-Linear Problem



第(1)

$x, x' \in X$

③ 神经网络. (XOR Prob)



XOR → 相同为0, 不同为1

$$x_1 \oplus x_2 = (\neg x_1 \wedge x_2) \vee (x_1 \wedge \neg x_2)$$

$$x_1 \text{ XOR } x_2 = \begin{cases} 1, 0 & \textcircled{1} \\ 0, 1 & \textcircled{2} \end{cases}$$

$$= \textcircled{1} \vee \textcircled{2}$$

$$= (x_1 \wedge \neg x_2) \vee (\neg x_1 \wedge x_2)$$

(OR: 包含)

$$1, 0 \rightarrow 1$$

$$0, 1 \rightarrow 1$$

$$1, 1 \rightarrow 1$$

(AND: 乘积)

$$1, 1 \rightarrow 1$$

(NOT)

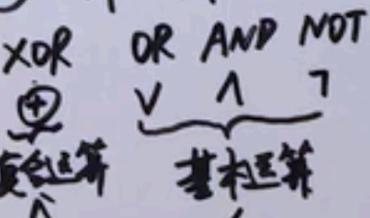
$$1 \rightarrow 0$$

第(2)

$K(x, x') = \langle \phi(x), \phi(x') \rangle$

$x, x' \in X$

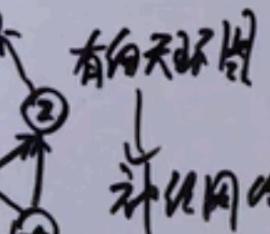
③ 神经网络. (XOR Prob)



XOR → 相同为0, 不同为1

$$x_1 \oplus x_2 = (\neg x_1 \wedge x_2) \vee (x_1 \wedge \neg x_2)$$

复合运算 → 复合表达式



$$\textcircled{1} \wedge$$

$$\textcircled{2} \wedge$$

$$\textcircled{3} \vee$$

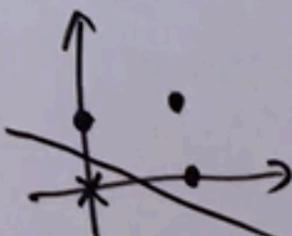
V

$$1, 0$$

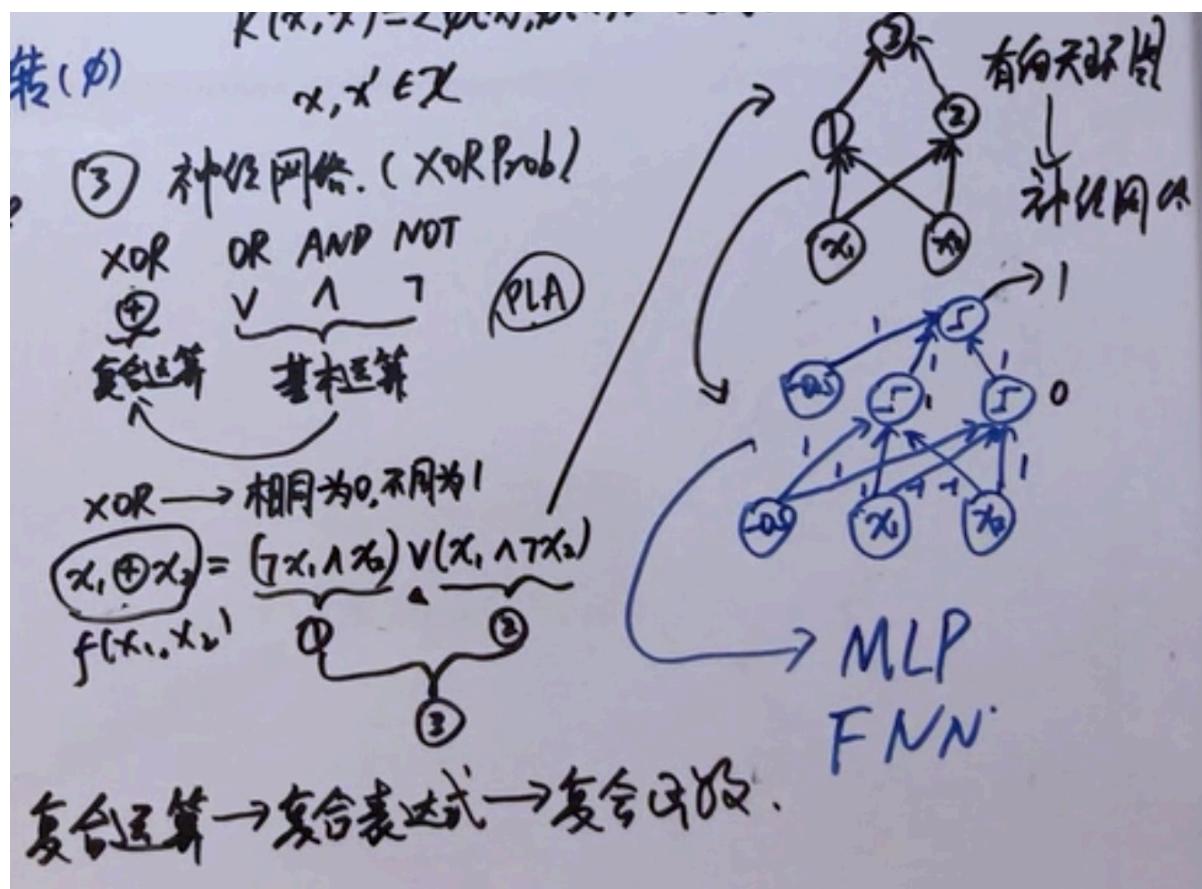
$$0, 1$$

$$1, 1$$

$$0, 0 \rightarrow 1$$



符合运算-》复合表达式-》复合函数



| 人工智能 两大阵营 三大主义

