

Logistic Regression as a Neural Network

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

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

Loss(error) function:

Don't use this, non-convex

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

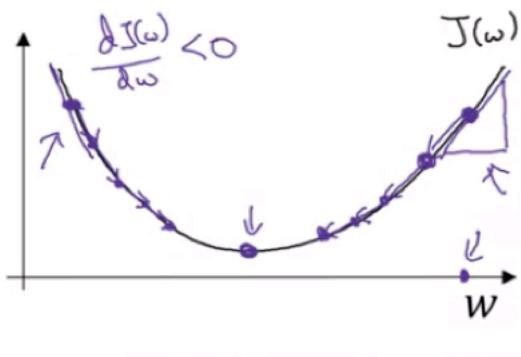
The loss function computes the error for a single training example; the cost function is the average of the loss functions of the entire training set.

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

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

a

Gradient Descent



$$J(w, b)$$

$$w := w - \alpha \frac{\partial J(w, b)}{\partial w}$$

$$b := b - \alpha \frac{\partial J(w, b)}{\partial b}$$

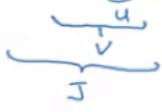
Repeat {
 $w := w - \alpha \frac{\partial J(w)}{\partial w}$
 $w := w - \alpha \frac{\partial J(w)}{\partial w}$
 $\frac{\partial J(w)}{\partial w} = ?$

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

$$\partial$$

Computation Graph

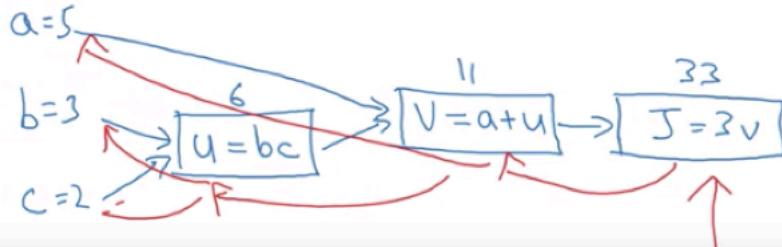
$$J(a, b, c) = 3(a + bc) = 3(5 + 3 \cdot 2) = 33$$



$$u = bc$$

$$v = a + u$$

$$J = 3v$$



Computing derivatives

Computing derivatives for $J(a, b, c) = 3(a + bc)$

Given values: $a = 5$, $b = 3$, $c = 2$

Intermediate values: $u = bc = 6$, $v = a + u = 11$, $J = 3v = 33$

Derivatives:

- $\frac{\partial J}{\partial a} = 3$
- $\frac{\partial J}{\partial b} = 3$
- $\frac{\partial J}{\partial c} = 9$
- $\frac{\partial J}{\partial u} = 3$
- $\frac{\partial J}{\partial b} = \frac{\partial J}{\partial u} \cdot \frac{\partial u}{\partial b} = 3 \cdot \frac{1}{2} = \frac{3}{2}$
- $\frac{\partial J}{\partial c} = \frac{\partial J}{\partial u} \cdot \frac{\partial u}{\partial c} = 3 \cdot 3 = 9$

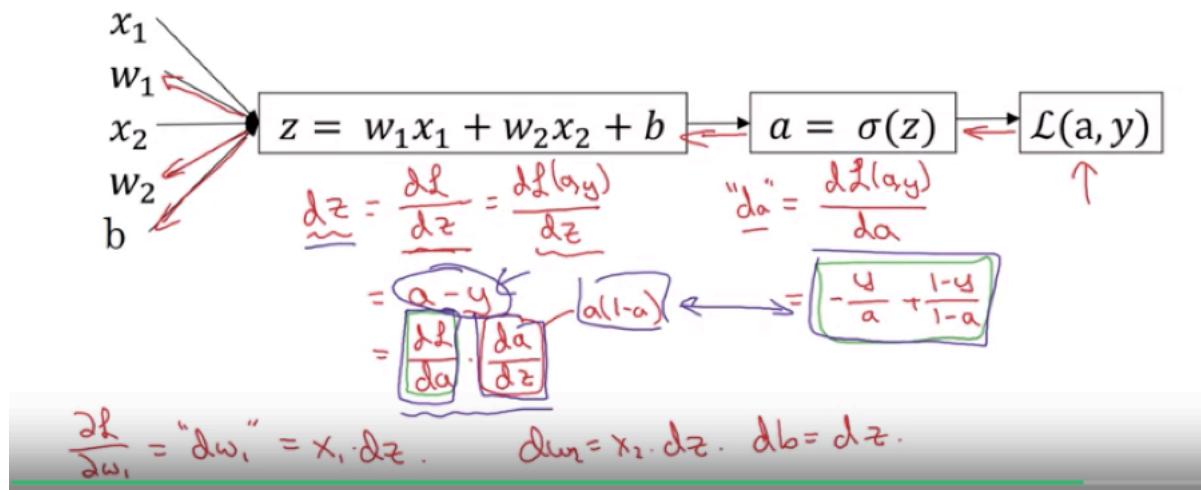
Final results:

- $u = 6 \rightarrow 6.001$
- $v = 11 \rightarrow 11.001$
- $J = 33 \rightarrow 33.003$
- $b = 3 \rightarrow 3.001$
- $u = b \cdot c = 6 \rightarrow 6.002$
- $J = 33.006$
- $v = 11.002$
- $J = 33.006$

Time: 13:32 / 14:33

Gradient Descent

Logistic regression derivatives



Logistic regression on m examples

$$\underline{J(\omega, b)} = \frac{1}{m} \sum_{i=1}^m \underline{\ell(a^{(i)}, y^{(i)})} \quad (\underline{x^{(i)}}, \underline{y^{(i)}})$$

$$\rightarrow \underline{a^{(i)}} = \underline{\hat{y}^{(i)}} = \sigma(\underline{z^{(i)}}) = \sigma(\underline{\omega^\top x^{(i)}} + b) \quad \underline{d\omega_1^{(i)}}, \underline{d\omega_2^{(i)}}, \underline{db^{(i)}}$$

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

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:

$$\begin{aligned}
\frac{dL}{dz} &= \frac{dL}{da} \times \frac{da}{dz} \\
\frac{dL}{dz} &= \frac{a - y}{a(1-a)} \times a(1 - a) = a - y
\end{aligned}$$

Vectorization

Logistic regression derivatives

$$\begin{aligned}
J &= 0, \quad \boxed{dw_1 = 0, dw_2 = 0}, \quad db = 0 \quad dw = np.zeros((n_x, 1)) \\
\rightarrow \text{for } i &= 1 \text{ to } m: \\
z^{(i)} &= w^T x^{(i)} + b \\
a^{(i)} &= \sigma(z^{(i)}) \\
J &+= -[y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})] \\
\frac{\partial J}{\partial w_j} &= x_j^{(i)} \frac{\partial}{\partial z^{(i)}} \left[a^{(i)}(1 - a^{(i)}) \right] \quad n_x = 2 \\
\frac{\partial J}{\partial w_1} &+= x_1^{(i)} dz^{(i)} \quad dw_1 += x_1^{(i)} dz^{(i)} \\
\frac{\partial J}{\partial w_2} &+= x_2^{(i)} dz^{(i)} \quad dw_2 += x_2^{(i)} dz^{(i)} \\
db &+= dz^{(i)} \quad db += dz^{(i)} \\
J &= J/m, \quad \boxed{dw_1 = dw_1/m, dw_2 = dw_2/m}, \quad db = db/m \\
dw &/= m
\end{aligned}$$

Vectorizing Logistic Regression

$$\begin{aligned}
 dz^{(1)} &= a^{(1)} - y^{(1)} & dz^{(2)} &= a^{(2)} - y^{(2)} & \dots \\
 dz &= [dz^{(1)} \ dz^{(2)} \ \dots \ dz^{(m)}]_m \leftarrow & & & \\
 A &= [a^{(1)} \ \dots \ a^{(m)}] & Y &= [y^{(1)} \ \dots \ y^{(m)}] \\
 \Rightarrow dz &= A - Y = [a^{(1)} - y^{(1)} \ a^{(2)} - y^{(2)} \ \dots]
 \end{aligned}$$

$$\begin{array}{l}
 \xrightarrow{\text{d}w_i} \begin{cases} \text{d}w = 0 \\ \text{d}w += \frac{x^{(1)} dz^{(1)}}{m} \\ \text{d}w += \frac{x^{(2)} dz^{(2)}}{m} \\ \vdots \\ \text{d}w / = m \end{cases} \\
 \xleftarrow{\text{d}z^{(i)}} \begin{cases} \text{d}b = 0 \\ \text{d}b += dz^{(1)} \\ \text{d}b += dz^{(2)} \\ \vdots \\ \text{d}b / = m \end{cases}
 \end{array}$$

$$\begin{aligned}
 db &= \frac{1}{m} \sum_{i=1}^m dz^{(i)} \\
 &= \frac{1}{m} \text{np.sum}(dz) \\
 dw &= \frac{1}{m} X dz^T \\
 &= \frac{1}{m} \left[\begin{matrix} x^{(1)} & \dots & x^{(m)} \end{matrix} \right] \left[\begin{matrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{matrix} \right] \\
 &= \frac{1}{m} \left[\underline{x^{(1)} dz^{(1)}} + \dots + \underline{x^{(m)} dz^{(m)}} \right]_{n \times 1}
 \end{aligned}$$

Implementing Logistic Regression

```

J = 0, dw1 = 0, dw2 = 0, db = 0
for i = 1 to m:
    z(i) = wTx(i) + b ←
    a(i) = σ(z(i)) ←
    J += -[y(i) log a(i) + (1 - y(i)) log(1 - a(i))]
    dz(i) = a(i) - y(i) ←
    dw1 += x1(i) dz(i)
    dw2 += x2(i) dz(i)
    db += dz(i)
J = J/m, dw1 = dw1/m, dw2 = dw2/m
db = db/m

```

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

Andrew

Logistic regression cost function

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

$$p(y|x) = \hat{y}^y (1-\hat{y})^{(1-y)} \quad \leftarrow$$

$$\text{If } y=1: p(y|x) = \hat{y} \stackrel{y}{=} (1-\hat{y})^0 = 1 \times (1-\hat{y}) = \underline{1-\hat{y}}$$

$$\text{If } y=0: p(y|x) = \hat{y}^0 \stackrel{0}{=} (1-\hat{y})^{(1-y)} = 1 \times (1-\hat{y}) = \underline{1-\hat{y}}$$

$$\uparrow \log p(y|x) = \log \hat{y}^y (1-\hat{y})^{(1-y)} = y \log \hat{y} + (1-y) \log (1-\hat{y})$$

$$= \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{(i)}, y^{(i)}) \downarrow$$

Andrew Ng

Cost on m examples

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

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
```

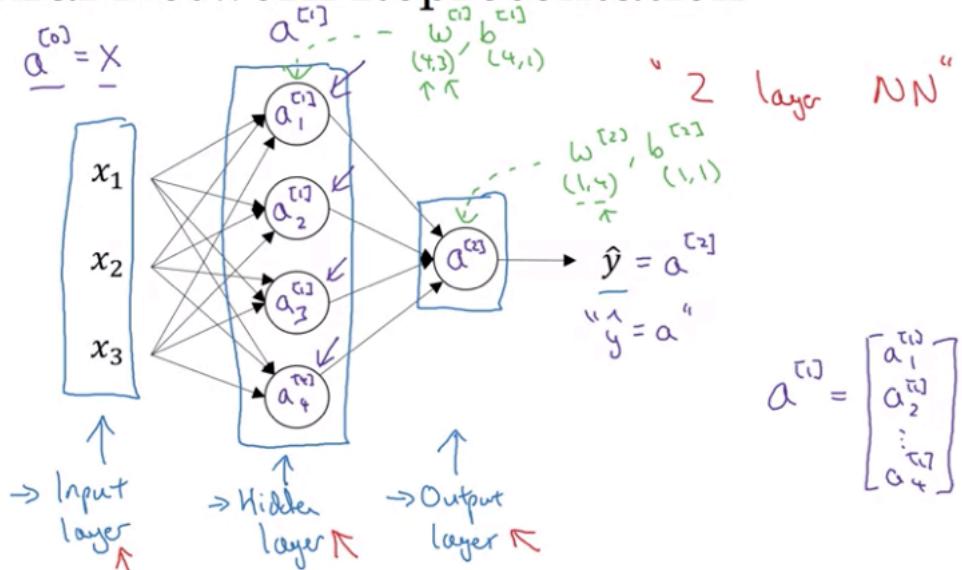
To represent color images, the red, green and blue channels (RGB) must be specified for each pixel, and so the pixel value is actually a vector of three numbers ranging from 0 to 255.

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).

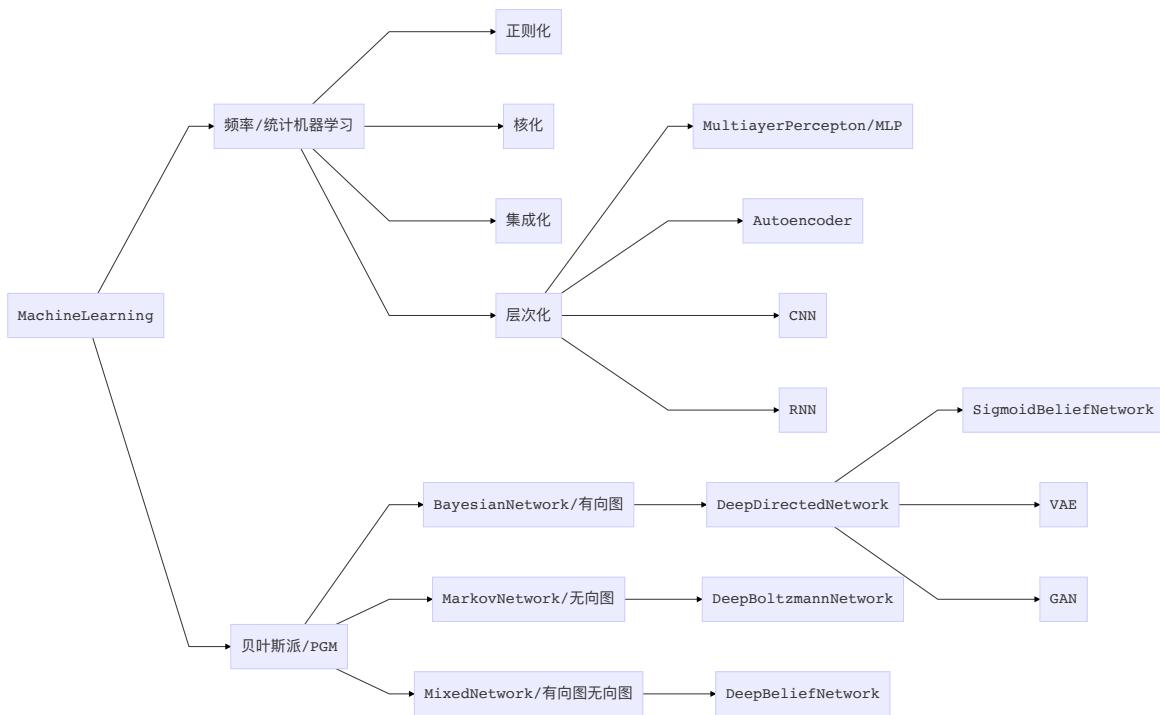
Let's standardize our dataset.

Neural Network

Neural Network Representation



基础



统计机器学习

- 正则化: Loss Function + regularizer
- 核化: Kernel SVM
- 集成化: AdaBoost, RandomForest
- 层次化: Neural Network/Deep Neural Network
 - MLP
 - 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

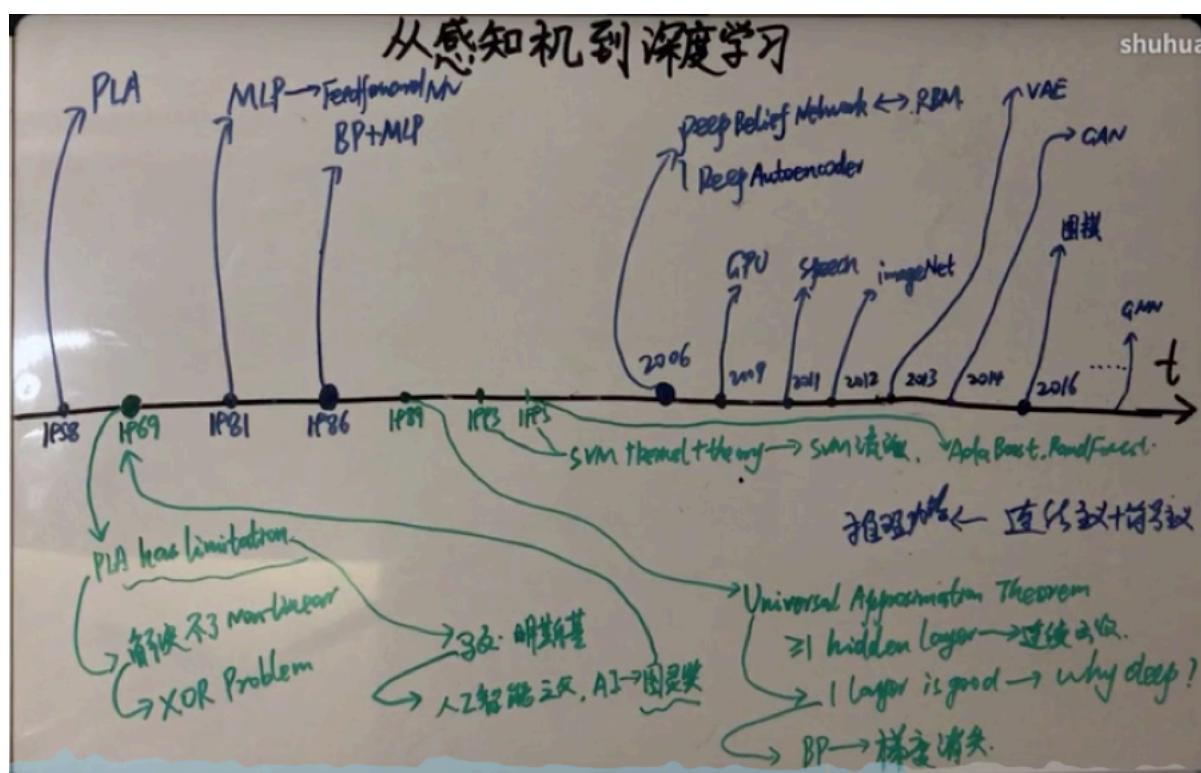
Deep Generative Model

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)

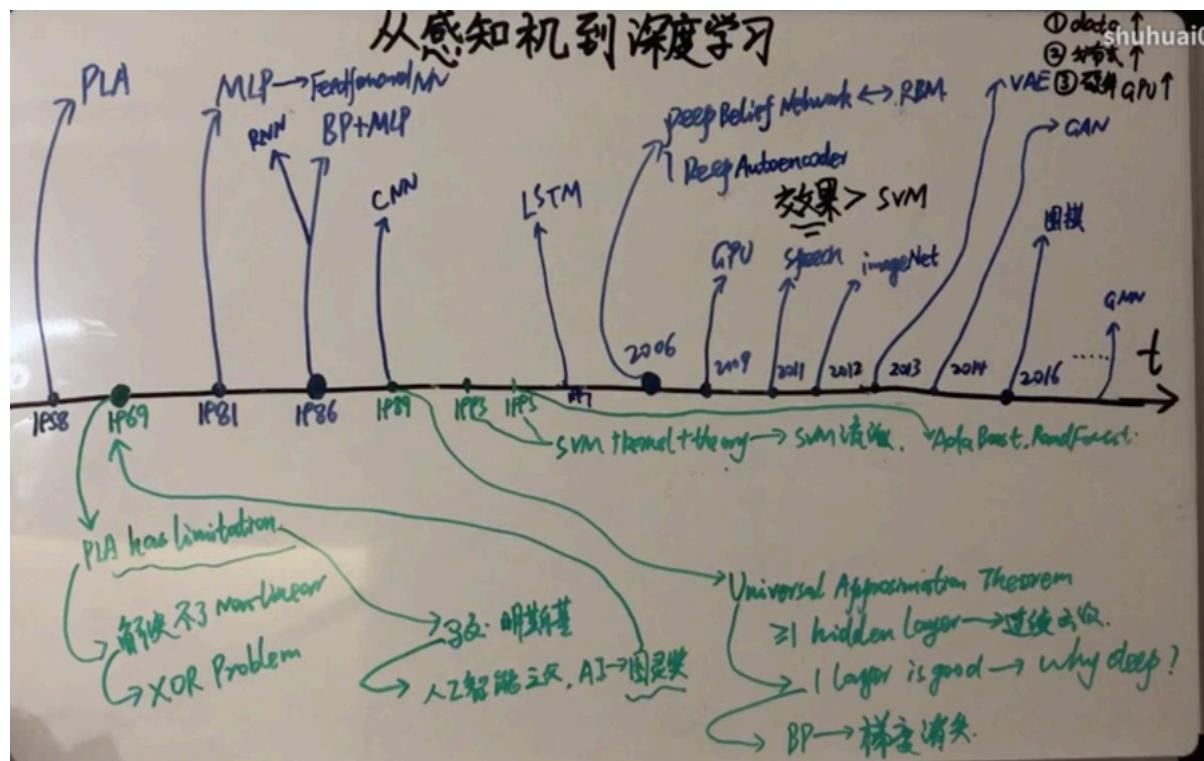
时间线



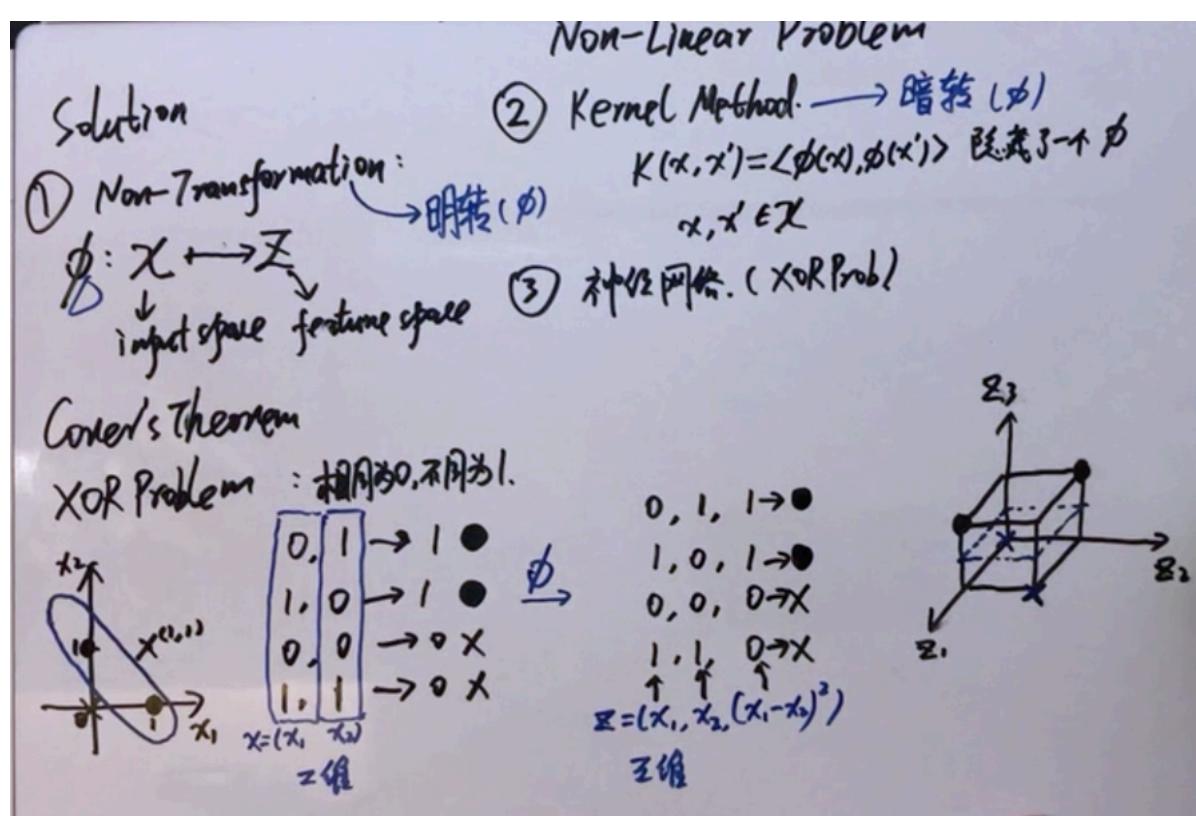
1. 深度学习的理论在2006年已经成型，直到现在，理论并没有根本性突破。

2. 为什么take off

1. data
2. 分布式
3. 硬件 GPU
4. 效果



Non-Linear Problem

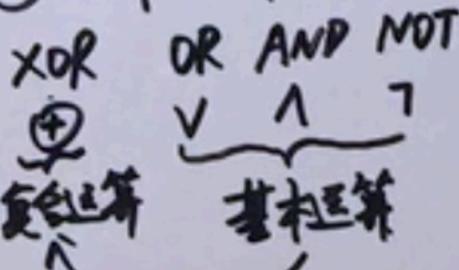


Neural Network

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

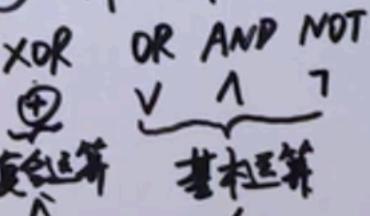
$$1 \rightarrow 0$$

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

OR: 包含

$$1, 0 \rightarrow 1$$

$$0, 1 \rightarrow 1$$

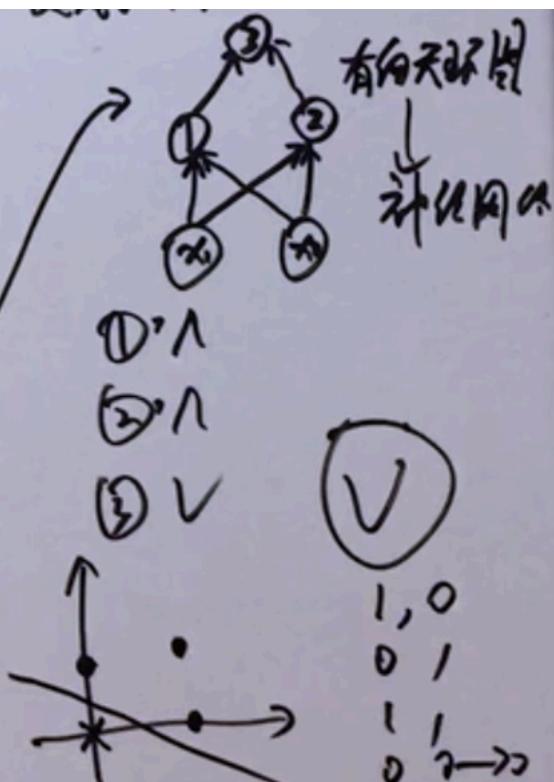
$$1, 1 \rightarrow 1$$

AND: 乘积

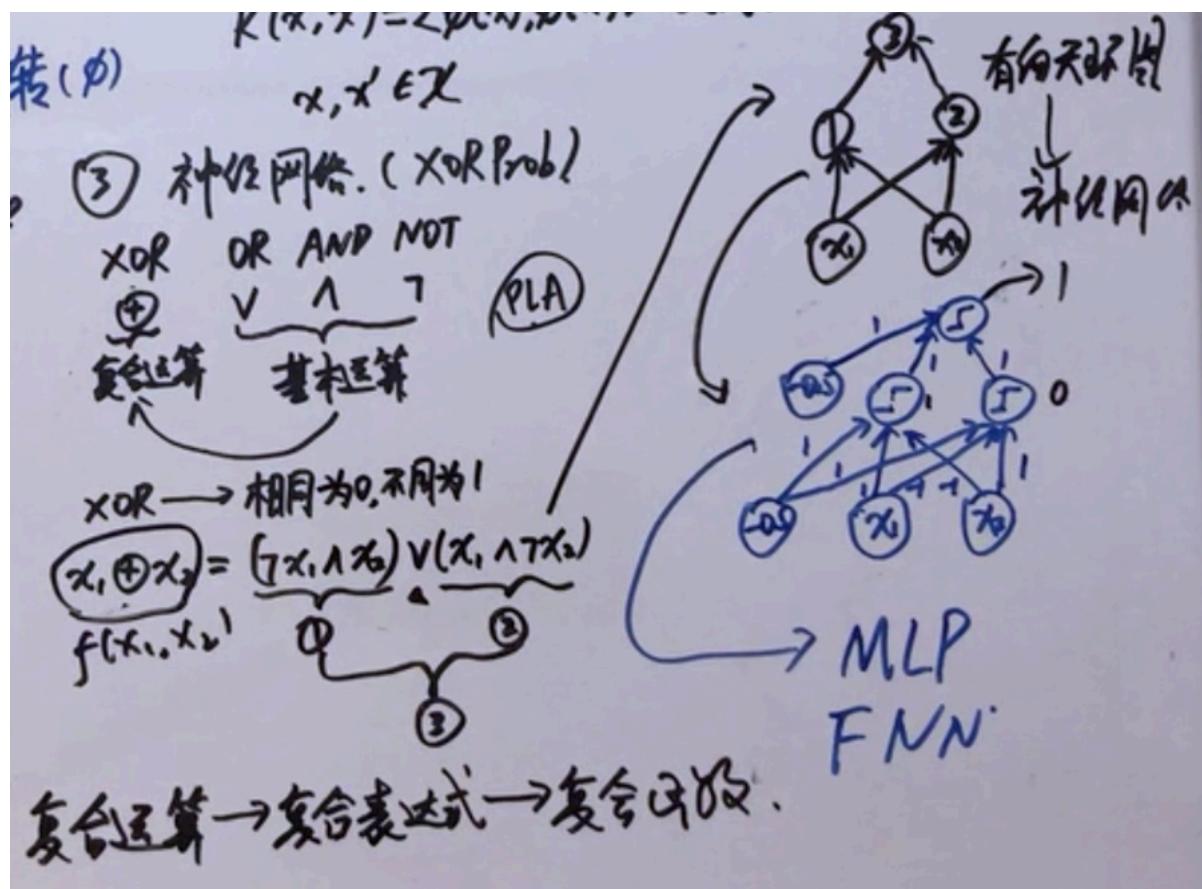
$$1, 1 \rightarrow 1$$

NOT

$$1 \rightarrow 0$$



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

