

Deep Learning (taught by Andrew Ng) – Notes by Xiaofan Wang

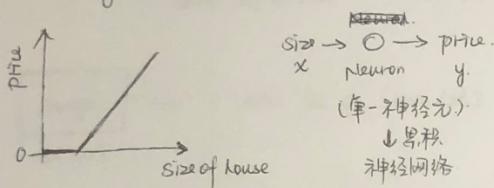
Name: Xiaofan Wang.

Neural Networks and Deep Learning.

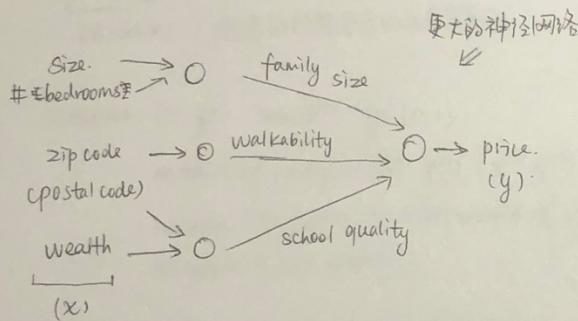
Ch1. Introduction to Deep Learning

1. Neural Network.

1. Housing Price Prediction.

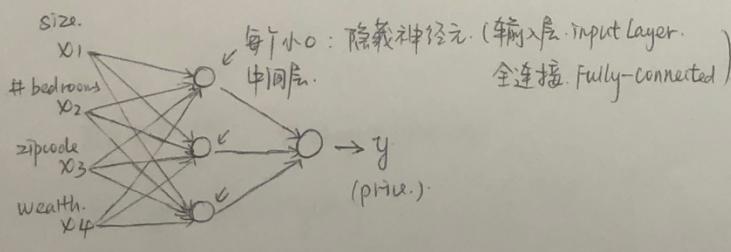


ReLU (Rectified Linear Unit. 线性整流函数)
max(0, y)



Neural Network: 只需给它训练集中的大量例子 $x \rightarrow y$.

而所有中间部分，NN会自己弄清。(只要给足够多的训练示例 $x \rightarrow y$)
NN就可以很好地拟合一个函数来建立 x 、 y 之间的映射关系)



输入参数: x_1, x_2, x_3, x_4 .

NN 的任务是: 预测 $\underline{\text{price}}$ "y"

2. Supervised Learning With Neural Network.

1.

Input (x)	Output (y)	Application.
① Home features	price.	Real Estate.
② Ad, User Info.	Click on Ad? 0/1.	Online Advertising.
③ Image.	Object (1---1000)	Photo tagging.
④ Audio	Text Transcript.	Speech Recognition.
⑤ English.	Chinese.	Machine Translation.
⑥ Image, Radar.	position of cars	Autonomous Driving.

在特定问题下，建立 $x \rightarrow y$ 映射关系，通过监督学习拟合数据，成为某个复杂系统的一部分。

①②: Standard NN.

③: CNN

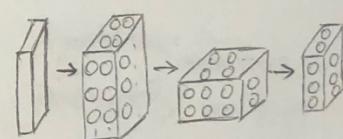
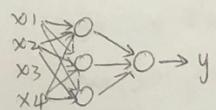
④⑥: RNN

⑤: Customized, hybrid NN architecture.

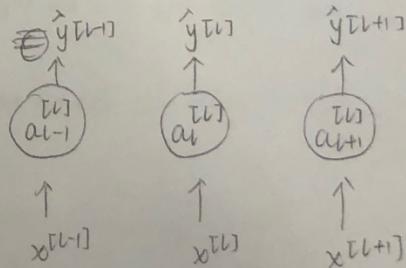
2. NN Examples:

Standard NN.

Convolutional NN.



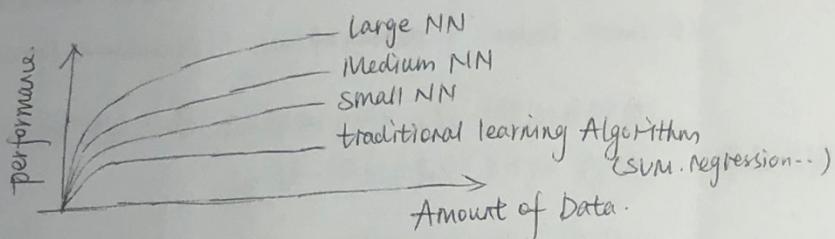
Recurrent NN. (循环神经网络)



3.
 { Structured Data: 每个特征都有清晰的含义
 { Unstructured Data:
 (Image, Text, Audio)

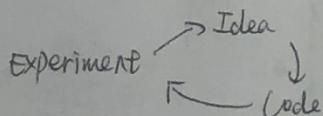
三. why is Deep Learning taking off?

1. Scale drives deep learning progress



High Performance: 足够大的网络. 大量数据

{ Data
 { Computation: sigmoid \rightarrow ReLU
 { Algorithms
 (\leftarrow) (\rightarrow)

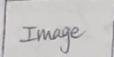


Ch2 Neural Network Basics

1. Logistic Regression as a Neural Network.

1. Binary Classification.

Neural Network 不使用 for 循环遍历全部训练样本。
使用“前向传播”、“反向传播”

 $\rightarrow 1 \text{ (cat)} \text{ vs } 0 \text{ (Non cat)}$
 64×64


 $\vec{x} = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$
 $n_x = 64 \times 64 \times 3 = 12288$
 $64 \times 64 \times 3$. (图像的矩阵包含的元素点数)

Notation: (x, y) . $x \in \mathbb{R}^{n_x}$, $y \in \{0, 1\}$.

m training example: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}) \dots (x^{(m)}, y^{(m)})\}$

M_{train} : emphasize the total number of sample.

$M_{\text{test}} = \# \text{ test example.}$

$$\vec{x} = \left[\begin{array}{c|c|c|c} & & & \vdots \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ \hline & & & \vdots \end{array} \right] \quad \underbrace{\quad}_{m} \quad \underbrace{n_x \quad \quad \quad}_{\perp}$$

$$\vec{x} \in \mathbb{R}^{n_x \times m} \quad \vec{x}. \text{shape} = (n_x, m)$$

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

$$\vec{y} \in \mathbb{R}^{l \times m}$$

$$\vec{y}. \text{shape} = (l, m)$$

2. Logistic Regression

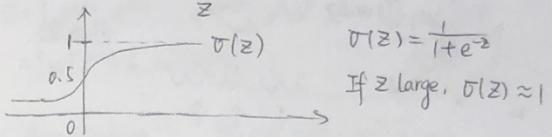
Given x (Image), $\Rightarrow \hat{y} = p(y=1|x) \in (0, 1)$.

输入图片 x , \hat{y} 告诉你这是-张猫图的概率。

$$\vec{x} \in \mathbb{R}^{n_x}$$

Parameters: $\vec{w} \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$.

$$\text{Output: } \hat{y} = \underline{\vec{w}^T \vec{x} + b}$$



$$x_0 = 1, \vec{x} \in \mathbb{R}^{n_x + 1}, \hat{y} = \sigma(\vec{w}^T \vec{x})$$

$$\vec{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_{n_x} \end{bmatrix} \quad \vec{b}$$

3. Logistic Regression Cost Function. (minimize it).

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

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

$$z^{(i)} = \vec{w}^T \vec{x}^{(i)} + b$$

$$\text{Loss (error) Function: } L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$$

$$L(\hat{y}, y) = -y \log \hat{y} + (1-y) \log(1-\hat{y})$$

If $y=1 \Rightarrow L(\hat{y}, y) = -\log \hat{y}$, want $\log \hat{y}$ large
want \hat{y} large.

If $y=0 \Rightarrow L(\hat{y}, y) = -\log(1-\hat{y})$, want $\log(1-\hat{y})$ large
want \hat{y} small.

$$\text{Cost Function: } J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

$$= -\frac{1}{m} \cdot \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)})]$$

寻找 w, b .

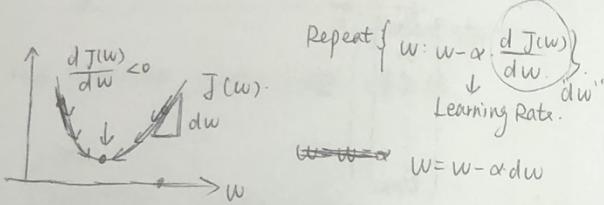
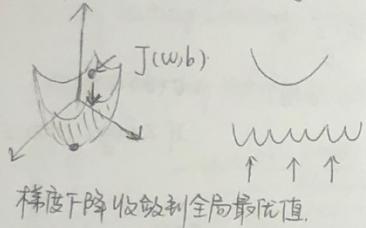
来最小化整体成本。

4 Gradient Descent. (学习训练集 $\Rightarrow w, b$)

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

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

want to find w, b that minimize $J(w, b)$.



$$J(w, b) : \begin{cases} w = w - \alpha \cdot \frac{dJ(w, b)}{dw} \\ b = b - \alpha \cdot \frac{dJ(w, b)}{db} \end{cases}$$

对 $\frac{dJ(w, b)}{dw}$: 当函数只有一个变量
 $\frac{\partial J(w, b)}{\partial w}$: 偏导数, 关其中一变量在对应点处的斜率.

6 Logistic Regression Gradient Descent.

① Logistic Regression Recap (单-样本).

$$z = w^T x + b$$

$$\hat{y} = a = \sigma(z)$$

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

$$\begin{aligned} x_1 &\rightarrow z = w_1 x_1 + w_2 x_2 + b \rightarrow a = \sigma(z) \rightarrow L(a, y) \\ x_2 &\rightarrow \\ b &\rightarrow \frac{d\sigma}{dz} = \frac{dL}{da} \quad "da" = \frac{dL(a, y)}{da} \\ &= \frac{dL(a, y)}{d\sigma} = -\frac{y}{a} + \frac{1-y}{1-a} \\ &= a - y. \end{aligned}$$

$$\begin{aligned} &= \frac{dL}{da} \left(\frac{da}{dz} \right) \rightarrow \text{all } a \\ &\left\{ \begin{array}{l} w_1 = w_1 - \alpha dw_1 \\ w_2 = w_2 - \alpha dw_2 \\ b = b - \alpha db \end{array} \right. \end{aligned}$$

② On M examples.

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(a^{(i)}, y^{(i)})$$

$$a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(w^T x^{(i)} + b)$$

$$\frac{\partial J(w, b)}{\partial w_1} = \frac{1}{m} \sum_{i=1}^m \frac{\partial L(a^{(i)}, y^{(i)})}{\partial w_1}$$

$$d w_1^{(i)} = (x^{(i)}, y^{(i)})$$

$$J = 0, dw_1 = 0, dw_2 = 0, db = 0.$$

For $i=1$ to m ,

$$z^{(i)} = w^T x^{(i)} + b$$

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

$$J += -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log (1-a^{(i)})]$$

(-VX下降的进阶) $d z^{(i)} = a^{(i)} - y^{(i)}$

$$\begin{aligned} \text{变量化} & \left\{ \begin{array}{l} dw_1 += x_1^{(i)} \cdot dz^{(i)} \\ dw_2 += x_2^{(i)} \cdot dz^{(i)} \\ db += dz^{(i)} \end{array} \right. \\ \text{避免偏杯} & \uparrow n=2 \\ & \downarrow \\ & \left\{ \begin{array}{l} dw_1 = w_1 - \alpha dw_1 \\ w_2 = w_2 - \alpha dw_2 \\ b = b - \alpha db \end{array} \right. \end{aligned}$$

$$J/m$$

$$dw_1/m, dw_2/m, db/m$$

$$[dw = np, zeros(n-x, 1)]$$

5 Computation Graph.

正向传播: 前向计算 NN Output.

反向传播: 计算梯度或微分

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

$$\begin{cases} u = bc \\ v = a+u \\ J = 3v \end{cases}$$

$$\begin{aligned} a &= 5 & \frac{dJ}{da} = da = 3 \\ b &= 3 & \frac{dJ}{db} = db = 3 \\ c &= 2 & \frac{dJ}{dc} = dc = 3 \\ u &= bc = 6 & \frac{du}{db} = du/db = 2 \\ & & \frac{du}{dc} = du/dc = 1 \\ & & \frac{du}{da} = du/da = 0 \end{aligned}$$

$$\begin{aligned} \frac{dJ}{dt} &= \frac{dJ}{du} \frac{du}{dt} \\ \frac{dJ}{dt} &= \frac{dJ}{db} \frac{db}{dt} + \frac{dJ}{dc} \frac{dc}{dt} \\ &= 3 \times 2 = 6. \end{aligned}$$

反向传播: $\frac{dJ}{dv} = 3$.

$$\frac{dJ}{dv} = 3 = \frac{dJ}{da} \frac{da}{dv}$$

$$= 3$$

$$= 1$$

$$= 1$$

$$= 1$$

$$= 1$$

$$= 1$$

$$= 1$$

$$= 1$$

$$= 1$$

$$= 1$$

$$= 1$$

$$= 1$$

$$= 1$$

$$= 1$$

2. Vectorization. (使代码艺术高效)

1. Vectorization.

$$z = w^T \cdot x + b$$

$$\vec{w} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \quad \vec{x} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

$\vec{w} \in \mathbb{R}^{n_x}$

Non-vectorized:

$$z = 0$$

for i in range(n_x):

$$z += w^{(i)} * x^{(i)}$$

$$z += b$$

Vectorized:

$$z = np.\text{dot}(\vec{w}, \vec{x}) + b$$

$$\vec{w}^T \cdot \vec{x}$$

3. Vectorizing Logistic Regression.

$$z^{(1)} = w^T \cdot x^{(1)} + b \quad z^{(2)} = w^T \cdot x^{(2)} + b \quad z^{(3)} = w^T \cdot x^{(3)} + b$$

$$a^{(1)} = \sigma(z^{(1)}) \quad a^{(2)} = \sigma(z^{(2)}) \quad a^{(3)} = \sigma(z^{(3)})$$

$$\vec{v} = \begin{bmatrix} \cdot & \cdot & \cdot \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ \cdot & \cdot & \cdot & \cdot \end{bmatrix} \quad (n_x, m)$$

$$\mathbb{R}^{n_x \times m}$$

$$z = [z^{(1)}, z^{(2)}, \dots, z^{(m)}] = w^T \cdot \vec{x} + [b, b, \dots, b]$$

$$= [w^T x^{(1)} + b, w^T x^{(2)} + b, \dots] \quad \underbrace{[b, b, \dots, b]}_{1 \times m}$$

$$(\text{"Broadcasting"}) \quad z = np.\text{dot}(w^T, \vec{x}) + \underbrace{b}_{(1, 1)}$$

$$\vec{A} = [a^{(1)}, a^{(2)}, \dots, a^{(m)}] = \sigma(z)$$

GPU \Rightarrow SIMD - single intuition multiple data.

CPU (通过内置函数和避免显示for循环，实现向量化。

提高运行速度)

2. More Examples.

① NN programming Guideline: whenever possible, avoid explicit for-loops.

Non-vectorized:

$$u = A \cdot v$$

$$u_i = \sum A_{ij} \cdot v_j$$

$$u = np.zeros(n, 1)$$

for i---

for j---

$$u[i] += A[i][j] * v[j]$$

Vectorized:

$$u = np.\text{dot}(A, v)$$

(消除两层for-loops
faster speed for operation)

4. Vectorizing LR's Gradient Computation.

$$\vec{A} = [a^{(1)}, \dots, a^{(m)}]$$

$$\vec{F} = [y^{(1)}, \dots, y^{(m)}]$$

$$dz = \vec{A} - \vec{F} = [a^{(1)}, y^{(1)}, \dots, a^{(m)}, y^{(m)}]$$

$$dw = 0$$

$$dw += x^{(1)} dz^{(1)}$$

:

$$dw / = m$$

$$db = 0$$

$$db += dz^{(1)}$$

:

$$db / = m$$

$$db = \frac{1}{m} \sum_{i=1}^m dz^{(i)}$$

$$= \frac{1}{m} \cdot np.sum(dz)$$

$$dw = \frac{1}{m} \cdot \vec{x} \cdot dz^T$$

$$= \frac{1}{m} \cdot \begin{bmatrix} x^{(1)}, x^{(m)} \end{bmatrix} \cdot \begin{bmatrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{bmatrix}$$

$$= \frac{1}{m} \cdot [x^{(1)} dz^{(1)}, \dots]$$

$$nx_1$$

LR's Gradient Computation - 使用法:

$$z = w^T \cdot x + b = np.\text{dot}(w^T, x) + b$$

$$\vec{A} = \sigma(z)$$

$$dz = \vec{A} - \vec{F}$$

$$dw = \frac{1}{m} \cdot \vec{x} \cdot dz^T$$

$$db = \frac{1}{m} \cdot np.sum(dz)$$

$$w = w - \alpha \cdot dw$$

$$b = b - \alpha \cdot db$$

② Vectors and matrix valued functions.

Need to apply the exponential operation on every element of a matrix/vector.

$$\vec{v} = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix} \quad \vec{u} = \begin{bmatrix} e^{v_1} \\ e^{v_2} \\ \vdots \\ e^{v_n} \end{bmatrix}$$

$$u = np.zeros(n, 1)$$

for i in range(n):

$$u[i] = \text{math.exp}(v[i])$$

Numpy 内置函数: import numpy as np
 $u = np.exp(v)$.

$$(np.log(v), np.abs(v), np.maximum(v, 0), v**2, v/v)$$

5 Broadcasting (faster, fewer codes)

Example:

Calories from Carbs, Proteins, Fats in 100g of different foods:

$$\begin{array}{l} \text{Carb: } \begin{array}{cccc} \text{Apples} & \text{Beef} & \text{Eggs} & \text{Potatoes} \\ [56.0 & 0.0 & 4.4 & 68.0] \end{array} \\ \text{Protein: } \begin{array}{cccc} 1.2 & 104.0 & 52.0 & 8.0 \\ [1.2 & 104.0 & 52.0 & 8.0] \end{array} = \vec{A}(3,4) \\ \text{Fat: } \begin{array}{cccc} 1.8 & 135.0 & 99.0 & 0.9 \\ [1.8 & 135.0 & 99.0 & 0.9] \end{array} \end{array}$$

Python. Broadcasting:

`cal = A.sum(axis=0)` # 计算每列和.

`percentage = 100 * A / (cal.reshape(1,4))`

(3,4)/(1,4) 当不确定矩阵维数时可用 reshape.

以保证是正确的行/列向量.

`reshape: O(1)`. 消耗成本低.

Other Example:

$$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + \begin{bmatrix} 100 \\ 100 \\ 100 \end{bmatrix} = \begin{bmatrix} 101 \\ 102 \\ 103 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \end{bmatrix} = \begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix}$$

General principle: $(m,n) + -*/ (1,n) \rightarrow (m,n)$
Matrix $(m,1) \rightarrow (m,n)$.

$(m,1) + R$.

$$[1, 2, 3] + 100 = [101, 102, 103]$$

$$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + 100 = \begin{bmatrix} 101 \\ 102 \\ 103 \end{bmatrix}$$

6. python./ Numpy Vectors (eliminate bug)

`a=np.random.randn(5). [a.shape=(5,1)].`

`a=np.random.randn(5,1) \Rightarrow a.shape=(5,1).` col vector

`a=np.random.randn(1,5) \Rightarrow a.shape=(1,5)` row vector

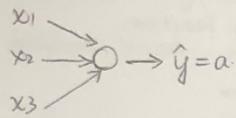
`assert(a.shape == (5,1)) \Rightarrow a=a.reshape((5,1)).`

用 assert 来审查矩阵和数组的维度.

Ch3 Shallow Neural Network

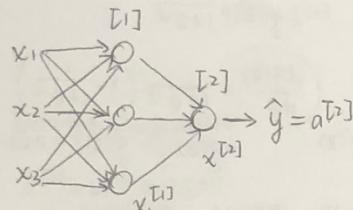
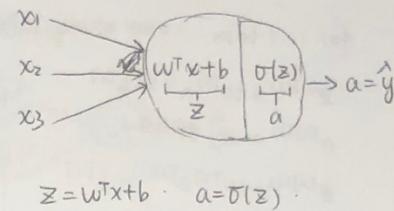
三. Computing NN's output.

一. Neural Network.



$$x \rightarrow z = w^T x + b \rightarrow a = \sigma(z) \rightarrow L(a, y)$$

$\underbrace{d z}_{da}$

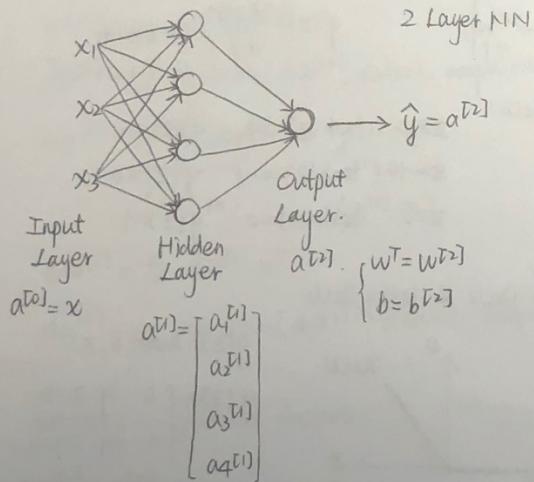


$$x^{[1]} \rightarrow z^{[1]} = w^{[1]} x + b^{[1]} \rightarrow a^{[1]} = \sigma(z^{[1]}) \rightarrow z^{[2]} = w^{[2]} \underbrace{a^{[1]} + b^{[2]}}_{a^{[2]}} \rightarrow a^{[2]} = \sigma(z^{[2]}) \rightarrow L(a^{[2]}, y)$$

$\underbrace{d z^{[2]}}_{d a^{[2]}}$

下指標^[1]的第一節點：
 $z^{[1]} = w_1^{[1] T} x + b_1^{[1]}$
 $a_1^{[1]} = \sigma(z_1^{[1]})$
 $a_i^{[1]} : \text{Node in layer.}$

二. Neural Network Representation.



parameters
in Hidden Layer. $\left\{ \begin{array}{l} w^{[1]} : (4, 3) \\ b^{[1]} : (4, 1) \end{array} \right.$

Hidden Layer. 2nd node:

$$z_2^{[1]} = w_2^{[1] T} x + b_2^{[1]}, a_2^{[1]} = \sigma(z_2^{[1]})$$

3rd Node; 4th node:

$$z_3^{[1]} = w_3^{[1] T} x + b_3^{[1]}, a_3^{[1]} = \sigma(z_3^{[1]})$$

$$z_4^{[1]} = w_4^{[1] T} x + b_4^{[1]}, a_4^{[1]} = \sigma(z_4^{[1]})$$

$$z^{[1]} = \begin{bmatrix} -w_1^{[1] T} & \dots \\ -w_2^{[1] T} & \dots \\ -w_3^{[1] T} & \dots \\ -w_4^{[1] T} & \dots \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{bmatrix} = \sigma(z^{[1]})$$

III. Vectorizing across multiple examples.

for $i=1 \text{ to } m$:

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

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

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

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

Given x :

$$x = A^{(0)} = a^{(0)i}(i)$$

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

$$A^{(1)} = \sigma(z^{(1)})$$

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

$$A^{(2)} = \sigma(z^{(2)})$$

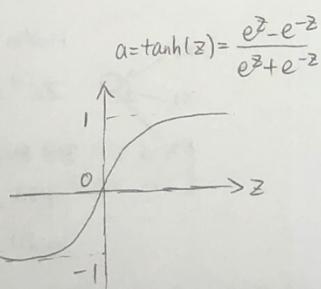
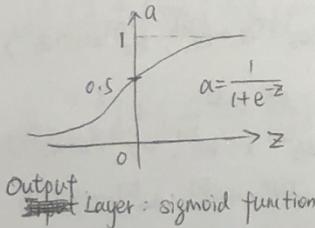
$$\vec{x} = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix}$$

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

$(n \times, m)$

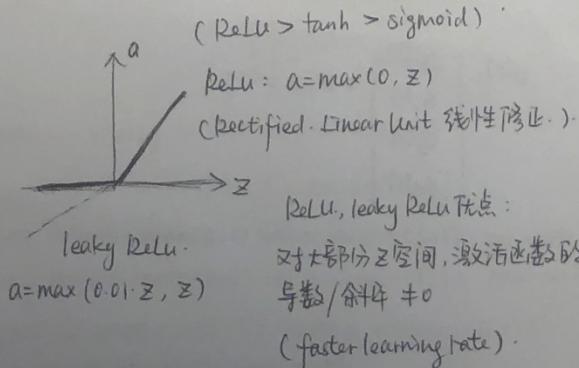
$$A^{(1)} = \begin{bmatrix} | & | & | \\ A^{(1)(1)} & A^{(1)(2)} & \dots & A^{(1)(m)} \\ | & | & | \end{bmatrix}$$

IV. Activation Functions



If z is very large/small,
 $k=a \approx 0$. decrease the speed
of gradient descent.

Hidden Layer: tanh function.
 $(\tanh > \text{sigmoid})$



V. Non-linear Activation Functions

deep NN: many hidden layers

t. Derivatives of Activation Functions

1. Sigmoid activation function



$$\frac{dg(z)}{dz} = \frac{1}{1 + e^{-z}} \left(1 - \frac{1}{1 + e^{-z}} \right)$$

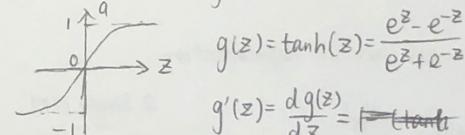
$$g'(z) = g(z) \cdot (1 - g(z)) = a(a-1)$$

Check: $z \text{ large } = 10$. $g(z) \approx 1$. $g'(z) \approx 0$

$z \text{ small } = -10$. $g(z) \approx 0$. $g'(z) \approx 0$

$z = 0$. $g(z) = \frac{1}{2}$. $g'(z) = \frac{1}{4}$

2. tanh activation function



$$g'(z) = \frac{d g(z)}{d z} = \frac{1 - \tanh^2(z)}{1 + \tanh^2(z)} = 1 - a^2$$

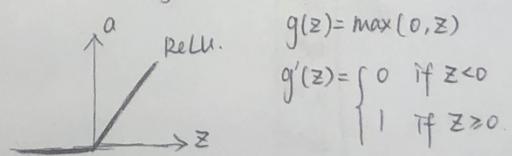
Check:

$z = 10$. $\tanh(z) \approx 1$. $g'(z) \approx 0$

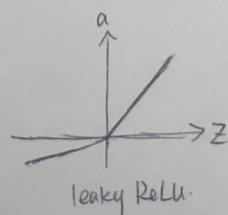
$z = -10$. $\tanh(z) \approx -1$. $g'(z) \approx 0$

$z = 0$. $\tanh(z) = 0$. $g'(z) = 1$.

3. ReLU & leaky ReLU



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



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

11. Gradient Descent for Neural Networks

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

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

Gradient Descent: $d w^{[1]} = \frac{\partial J}{\partial w^{[1]}} \cdot d b^{[1]} = \frac{\partial J}{\partial b^{[1]}} \dots$

$$w^{[1]} = w^{[1]} - \alpha \cdot d w^{[1]}$$

$$b^{[1]} = b^{[1]} - \alpha \cdot d b^{[1]}$$

.....

Forward Propagation:

$$z^{[1]} = w^{[1]}x + b^{[1]}$$

$$a^{[1]} = g^{[1]}(z^{[1]})$$

$$z^{[2]} = w^{[2]}A^{[1]} + b^{[2]}$$

$$A^{[2]} = g^{[2]}(z^{[2]}) = \sigma(z^{[2]})$$

Backward propagation:

$$d z^{[2]} = A^{[2]} - y \quad F = [y^{(1)}, y^{(2)} \dots y^{(m)}]$$

$$d w^{[2]} = \frac{1}{m} \cdot d z^{[2]} A^{[1]T} \quad (n^{[2]}, 1)$$

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

$$d z^{[1]} = \underbrace{w^{[2]T} \cdot d z^{[2]} * g'^{[2]}(z^{[1]})}_{(n^{[1]}, m)} \quad (n^{[1]}, m)$$

$$d w^{[1]} = \frac{1}{m} \cdot d z^{[1]} x^T$$

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

Summary of gradient descent:

$$d z^{[2]} = A^{[2]} - y$$

$$d w^{[2]} = \frac{1}{m} \cdot d z^{[2]} A^{[1]T}$$

$$d b^{[2]} = d z^{[2]}$$

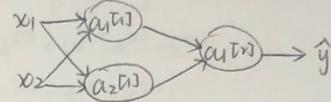
$$d z^{[1]} = w^{[2]T} \cdot d z^{[2]} * g'^{[2]}(z^{[1]})$$

$$d w^{[1]} = d z^{[1]} x^T$$

$$d b^{[1]} = d z^{[1]}$$

fr. Random Initialization.

If initialize weights to zero:



$$w^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}, \quad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad w^{[2]} = \begin{bmatrix} 0 & 0 \end{bmatrix}$$

$$a_1^{[1]} = a_2^{[1]} \quad d z_1^{[1]} = d z_2^{[1]}$$

$$d w = \begin{bmatrix} u & v \\ u & v \end{bmatrix} \quad w^{[1]} = w^{[1]} - \alpha d w$$

$$w^{[1]} = \text{np.random.randn}(2, 2) * 0.01$$

$$b^{[1]} = \text{np.zeros}(2, 1)$$

$$w^{[2]} = \text{np.random.randn}(1, 2) * 0.01$$

$$b^{[2]} = 0$$

Tips: 如何初始化神经网络的权重
⇒ 随机初始化参数对训练神经网络很重要。

$$d z^{[2]} = A^{[2]} - y$$

$$d w^{[2]} = \frac{1}{m} \cdot d z^{[2]} A^{[1]T}$$

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

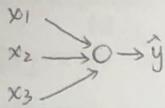
$$d z^{[1]} = w^{[2]T} \cdot d z^{[2]} * g'^{[2]}(z^{[1]})$$

$$d w^{[1]} = \frac{1}{m} \cdot d z^{[1]} x^T$$

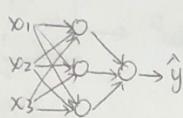
$$d b^{[1]} = \frac{1}{m} \cdot \text{np.sum}(d z^{[1]}, \text{axis}=1, \text{keepdims=True})$$

Ch 4. Deep Neural Network.

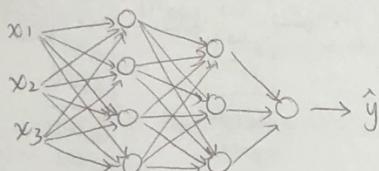
- Deep L-layer Neural Network.



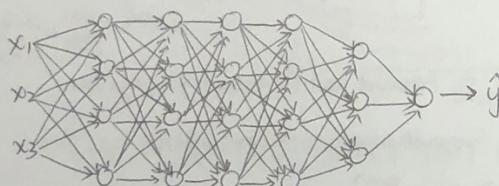
1 layer NN ("Shallow")
Logistic Regression.



1 Hidden Layer.
2 Layer NN.



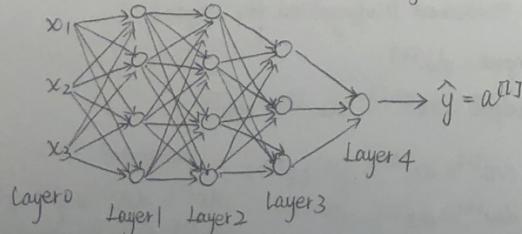
2 Hidden Layer.



5 Hidden Layer ("Deep").

Deep Neural Network Notation:

"4 Layer NN"



$n^{[l]}$: # units in layer l ($n^{[0]}=3$, $n^{[1]}=4$, $n^{[2]}=4$, $n^{[3]}=3$, $n^{[4]}=1$).

$a^{[l]}$: activation in layer l ($n^{[0]}=n_x=3$)

= Forward Propagation in a deep Network.

$$z^{[l]} = w^{[l]}x + b^{[l]} \quad (x = a^{[0]})$$

$$a^{[l]} = g^{[l]}(z^{[l]})$$

$$z^{[l+1]} = w^{[l+1]}a^{[l]} + b^{[l+1]}$$

$$a^{[l+1]} = g^{[l+1]}(z^{[l+1]})$$

⋮

$$z^{[L]} = w^{[L]}a^{[L-1]} + b^{[L]} \Rightarrow \begin{cases} z^{[l]} = w^{[l]}a^{[l-1]} + b^{[l]} \\ a^{[l]} = g^{[l]}(z^{[l]}) \end{cases}$$

"a", "z": dimensions: $(n^{[l]}, 1)$

Vectorized:

$$z^{[l]} = w^{[l]}A^{[0]} + b^{[l]} \quad (x = A^{[0]})$$

$$A^{[l]} = g^{[l]}(z^{[l]})$$

$$z^{[l+1]} = w^{[l+1]}A^{[l]} + b^{[l+1]}$$

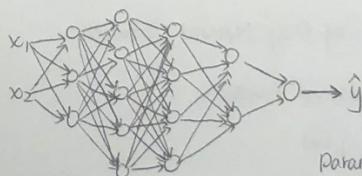
$$A^{[l+1]} = g^{[l+1]}(z^{[l+1]})$$

$$\hat{y} = g(z^{[L]}) = A^{[L]}$$

Deep Network \approx Network only with 1 Hidden Layer.

(只是多重复了几次).

= Get Matrix Dimensions Right



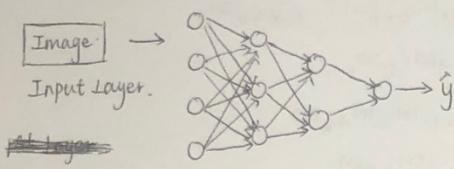
Parameters $w^{[l]}, b^{[l]}$

$$\begin{cases} w^{[l]}: (n^{[l]}, n^{[l-1]}) \\ b^{[l]}: (n^{[l]}, 1) \\ dw^{[l]}: (n^{[l]}, n^{[l-1]}) \\ db^{[l]}: (n^{[l]}, 1) \end{cases}$$

Vectorized Implementation:

$$\begin{cases} z^{[l]} = w^{[l]} \cdot \vec{x} + b^{[l]} : (n^{[l]}, 1) \\ z^{[l]}, A^{[l]}: (n^{[l]}, m) \\ dz^{[l]}, dA^{[l]}: (n^{[l]}, m) \\ l=0, A^{[0]} = \vec{x} : (n^{[0]}, m) \end{cases}$$

IV. Deep Representation.



- 1st Layer: Feature Detector / Edge Detector.
- 2nd Layer: combine "Edges" together \Rightarrow form parts of Image
Different Neuron \rightarrow Different part. (complex).
- 3rd Layer: combine "Different parts" into a whole.
 \Rightarrow Test Different Types of Faces.

先找寻“简单事物”

然后构建“简单事物”，将其组合在一起以检测“更复杂事物”
 \Downarrow

Detect "Simple Objects" $\xrightarrow{\text{combine}}$ "More Complex Objects"

Circuit theory and deep learning:

There are functions you can compute with a "small" L-layer deep neural network that shallower networks require exponentially more hidden units to compute.

V. Building Blocks of Deep Neural Network.

Forward and Backward Functions:

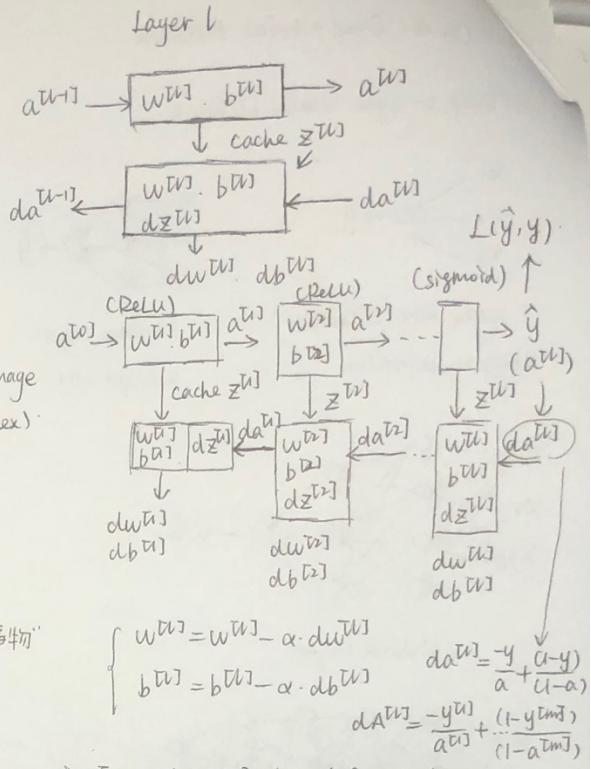
Layer l: $w^{[l]}, b^{[l]}$

Forward: Input $a^{[l-1]}$. Output $a^{[l]}$, Cache $z^{[l]}$

$$z^{[l]} = w^{[l]} \cdot a^{[l-1]} + b^{[l]}$$

$$a^{[l]} = g^{[l]}(z^{[l]})$$

Backward: Input $d_a^{[l]}$. Output $d_a^{[l-1]}, d_w^{[l]}, d_b^{[l]}$
Cache $z^{[l]}$.



vi. Forward and Backward Propagation.

1. Forward propagation for layer l.

Input: $a^{[l-1]}$

Output: $a^{[l]}$, Cache($z^{[l]}$) $\xrightarrow{w^{[l]}, b^{[l]}}$

$$\begin{cases} z^{[l]} = w^{[l]} \cdot a^{[l-1]} + b^{[l]} \\ a^{[l]} = g^{[l]}(z^{[l]}) \end{cases} \quad \begin{matrix} \text{Vectorized:} \\ z^{[l]} = w^{[l]}.A^{[l-1]} + b^{[l]} \\ A^{[l]} = g^{[l]}.(z^{[l]}) \end{matrix}$$

2. Backward propagation for Layer l.

Input: $d_a^{[l]}$

Output: $d_a^{[l-1]}, d_w^{[l]}, d_b^{[l]}$

$$\begin{cases} dz^{[l]} = d_a^{[l]} * g'^{[l]}(z^{[l]}) \\ dw^{[l]} = dz^{[l]}.a^{[l-1]} \\ db^{[l]} = dz^{[l]} \\ da^{[l-1]} = w^{[l]T}.dz^{[l]} \end{cases} \quad \begin{matrix} \text{Vectorized:} \\ dz^{[l]} = dA^{[l]} * g'^{[l]}(z^{[l]}) \\ dw^{[l]} = \frac{1}{m} \cdot dz^{[l]}.A^{[l-1]T} \\ db^{[l]} = \frac{1}{m} \cdot np.sum(dz^{[l]}), axis=1, keepdims=True \\ da^{[l-1]} = w^{[l]T}.dz^{[l]} \end{matrix}$$

(算法复杂度可能来源于数据而非代码)

c. Parameters. VS. Hyperparameters

(使运行高效，更好地配置超参数)

1. Hyperparameters

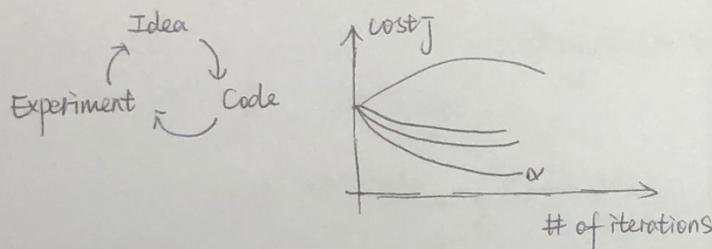
parameters: $w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)} \dots$

Hyperparameters: learning rate α .

(影响训练速度)
参数 w, b)

$\left. \begin{array}{l} \# \text{ iterations} \\ \# \text{ Hidden Layers } L \\ \# \text{ Hidden Units } (n^{(1)}, n^{(2)} \dots) \\ \text{Choice of Activation Functions.} \end{array} \right\}$

2. Applied Deep Learning is a very empirical process.



尝试许多不同的配置，构建模型并运行，看比较结果。

Name: Xiaofan Wang

Improving Deep Neural Network: Hyperparameter tuning, Regularization and Optimization.

Ch 1. Practical Aspects of Deep Learning

1. Set up Machine Learning Application.

1. Train / Dev / Test Sets.

Applied ML is a highly iterative process:

layers
hidden units
learning rates
activation functions

尝试找到一个越来越好的神经网络而不断迭代。

恰当地将 Dataset 分为 Train / Dev / Test sets 可提高迭代效率
(训练集 / 开发集 / 测试集)

Data	Training Set	Hold-out (Dev Set)	Test Set
训练算法	cross Validation	评估结果	
最佳	60% / 20% / 20%	测试许多模型	
tt	98% / 1% / 1%	比较哪个效果好	
BY	99.5% / 0.4% / 0.1%		

Mismatched Train / Test Distribution:

Training Set: Cat pictures from webpages (高清)
Dev / Test Sets: Cat pictures from users using your app (模糊)

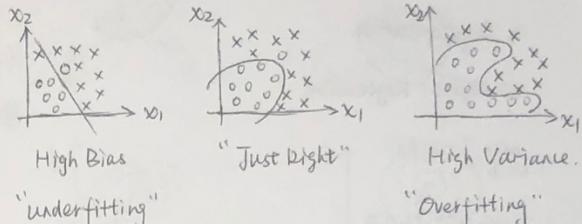
⇒ Make sure Dev / Test come from same distribution.

⇒ Not Having a Test Set might be Okay. (Only Dev Set)

测试集的目的是得到一个无偏估计 来评价最终选取的网络的性能。当只有开发集时，用训练集尝试不同的模型结构，并尝试得到一个好模型。(此时无需无偏估计模型)

2. Bias 偏差 (偏差度) & Variance 方差 (集中度).

"偏差一方差困境"



cat Classification: $\begin{cases} y=1 \\ y=0 \end{cases}$

Training Set Error: 1% | 15% | 15% | 0.5%

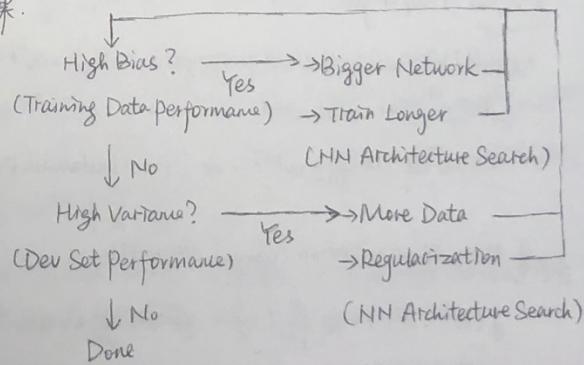
Dev Set Error: 11% | 16% | 30% | 1%

High Var. | High Bias | High Var. low V | High Bias low B

此时, 理想误差(贝叶斯误差) ≈ 0 .

Training / Dev Set from the same Distribution.

3. Basic recipe for Machine Learning.



Bias. Variance
↓ ↑
↑ ↓
(cost: computation time)
Bigger Network: Decrease Bias, No impact on Var
Regularization: Decrease Var, No impact on Bias.

2. Regularizing Neural Network

(防止高方差，过拟合)

1. Regularization

① Logistic Regression

$$\min_{w, b} J(w, b)$$

$w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|_2^2 + \frac{\lambda}{2m} b^2$$

$$L_2 \text{ regularization: } \|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w \quad \text{Omit.}$$

(w will be sparse)

$$\frac{\lambda}{2m} \cdot \sum_{j=1}^{n_x} |w_j| = \frac{\lambda}{2m} \|w\|_1$$

λ = regularization parameter

② Neural Network

$$J(w^{[0]}, b^{[0]}, \dots, w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_F^2$$

$$\|w^{[l]}\|_F^2 = \sum_{j=1}^{n^{[l-1]}} \sum_{j=1}^{n^{[l]}} (w_j^{[l]})^2, \quad w = (w^{[0]}, \dots, w^{[L-1]})$$

$$d[w^{[l]}] = (\text{from backward propagation}) + \frac{\lambda}{m} \cdot w^{[l]} \quad (\frac{\partial J}{\partial w^{[l]}} = d[w^{[l]}])$$

$$w^{[l]} = w^{[l]} - \alpha \cdot d[w^{[l]}]$$

$$= (1 - \alpha \cdot \frac{\lambda}{m}) w^{[l]} - \alpha \cdot (\text{from backward propagation}).$$

2. Why Regularization reduces Overfitting.

$$J(w^{[0]}, b^{[0]}, \dots, w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_F^2$$

$$w^{[l]} \approx 0$$

隐藏单元的权重接近0，隐藏单元的影响被消除了。

简化了Network，但很可能深度变大。~~所以~~（这个过拟合网络可能是高偏差状态，但可能存在一个中间值使 $w \approx 0$ ）

所以

不容易过拟合的小型神经网络。

3. Dropout regularization

先遍历 Network 的每层，每个节点的丢弃概率 0.5。

按随机编码决定每个节点的去留。

Implement Dropout:

layer = 3. keep_prob = 0.8. (随机设置使 a_3 不变)

```
d3 = np.random.rand(a3.shape[0], a3.shape[1])
a3 = np.multiply(a3, d3)
<keep_prob>
```

→ $a_3 / \text{keep_prob}$

50 units → 10 units shut off.

$$z^{[4]} = w^{[4]} \cdot a^{[3]} + b^{[4]}$$

↑
reduced by 20%.

$$/ = 0.8.$$

Making Predictions at Test Time:

$$a^{[0]} = \bar{x}$$

$$\text{No Dropout: } z^{[1]} = w^{[1]} a^{[0]} + b^{[1]}$$

$$a^{[1]} = g^{[1]}(z^{[1]})$$

$$z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

$$a^{[2]} = \dots$$

↓

$$\hat{y}$$

① 此时，不用抽硬币来决定哪些隐藏单元要被消除，不想让输出也是随机的，只会为预测增加噪声。

② 用不同的随机失活的 Network 进行多次预测，取平均值。⇒ 运算效率不高。

会得到几乎相同的预测~

③ 反向随机失活：保证如果测试过程没有对随机失活算法进行缩放 (scaling)，那么激活函数的期望输出也不会改变。所以，不用加入额外的缩放参数。

4. Understand Dropout.

Why does drop-out work?

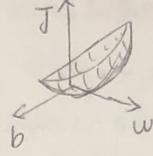
Intuition: Can't rely on any one feature.

So have to spread out weights. (shrink weights)

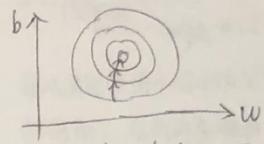
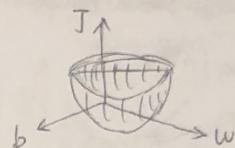
不要把太多的权重放在某个上. 对于每个输入都给一个较小权重.

Why normalize Inputs? $J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i)$

Unnormalized:



Normalized:



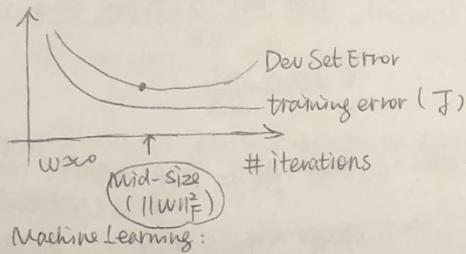
在梯度下降中采用更长的步长.

使特征尺度更近. 优化算法速度.

I. Other Regularization Methods.

① Data Augmentation. 数据扩增.

② Early Stopping. 早终止法. (提前终止训练)



Machine Learning:

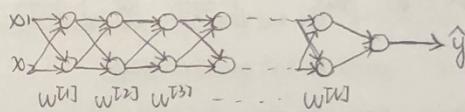
- Algorithm & Optimize "J" → Gradient Momentum RMSProp Adam
- Not Overfit: Regularization

Orthogonalization 正交化: 同一时间只考虑一个任务.

Early Stopping 优点: 只运行一次梯度下降的过程.

(可用L2代替) 要尝试小w值. 中w值. 大w值.
而不用尝试L2正则化中各种插入.

2. Vanishing / Exploding Gradients



$$g(z) = z, b^{(l)} = 0$$

$$\hat{y} = w^{(L)} \cdot w^{(L-1)} \cdot w^{(L-2)} \cdots w^{(3)} \cdot w^{(2)} \cdot w^{(1)} \cdot x$$

$$w^{(1)} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$$

$$z^{(1)} = w^{(1)} \cdot x$$

$$\hat{y} = w^{(L)} \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{-1} \cdot x$$

$$a^{(1)} = g(z^{(1)}) = z^{(1)}$$

$$a^{(L)} = g(z^{(L)}) = g(z^{(L)}, a^{(L)})$$

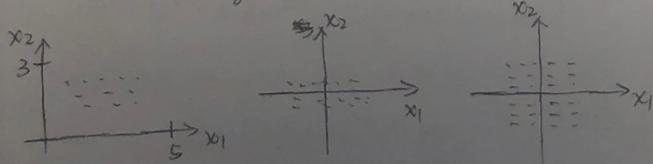
① 对于很深的 Network, L会很大. 以1.5的层数以方的速度变大. y的值就会爆炸

② 如果以0.5代替1.5. 激活函数的值会指数级下降. 此时梯度下降步长很小. 用很长时间才会产生学习.

III. Optimization.

1. Normalize Inputs.

Normalize Training Sets.



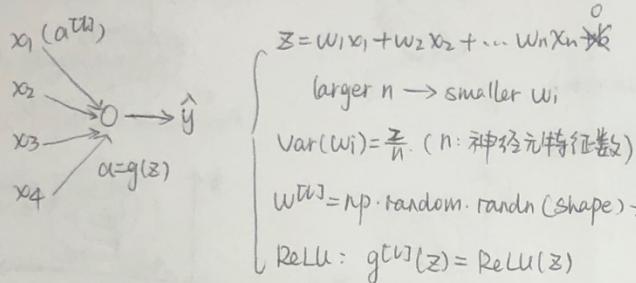
$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2$$

$$x - \mu$$

3. Weight Initialization for Deep Networks.

① Single Neuron Example: (激活函数是 ReLU)

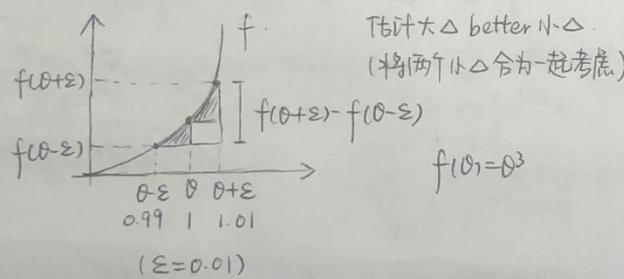


如果输入经过激活的特征平均值为0，且标准差为1。
就可使头顶也呈现相同的分布性质。
虽然不能完全解决问题，但通过设置权重 w 使 ~~w不会比~~
~~很多~~，使 $w \approx 1$ ($w \neq 1$)。因此梯度不会过快膨胀/消失。

② 激活函数: tanh. $\sqrt{\frac{1}{n^{[L-1]}}}$ 或 $\sqrt{\frac{2}{n^{[L-1]} + n^{[L]}}}$

4. Numerical Approximation of Gradients.

Check Derivative Computation:



$$g(\theta) \approx \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} = \frac{(1.01)^3 - (0.99)^3}{2 \times 0.01} = 3.0001$$

$$g(\theta) = 3\theta^2 = 3.$$

approximate error: $0.0001 \cdot (\text{双侧差值})$.

单侧差值: $3.0301 - 3 = 0.0301$.

所以，双侧差值 Better Than 单侧。

5. Gradient Checking.

① Take $w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}$ and reshape into a big vector θ .

$$\Rightarrow J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = J(\theta)$$

② Take $d\theta^{[1]}, db^{[1]}, \dots, d\theta^{[L]}, db^{[L]}$ and reshape into a big vector $d\theta$.

$$d\theta^{[1]} = np.\text{random}.\text{randn}(\text{shape}) * np.\text{sqrt}(\frac{2}{n^{[1]}})$$

$$J(\theta) = J(\theta_1, \theta_2, \dots, \theta_L)$$

for each i :

$$d\theta_{\text{approx}}^{[i]} = \frac{J(\theta_1 - \epsilon, \theta_2, \dots, \theta_L) - J(\theta_1, \theta_2, \dots, \theta_{i-1}, \theta_i + \epsilon, \theta_{i+1}, \dots)}{2\epsilon}$$

$$\approx d\theta_i = \frac{\partial J}{\partial \theta_i}$$

Check: $\|d\theta_{\text{app}} - d\theta\|_2 \approx 10^{-7} \rightarrow \text{great! (Debug)}$.
(10^{-5} to 10^{-3} , worry).

$$\Rightarrow \|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2 \quad \epsilon = 10^{-7}$$

Notes:

- Don't use in training. Only to debug.

$d\theta_{\text{approx}}^{[i]}$: slow., 当用梯度下降时，用反向传播计算 $d\theta$.

- If algorithm fails grad check, look at components to try to identify bug.

当 $d\theta_{\text{approx}}$ 与 $d\theta$ 差距很大. \Rightarrow 基层 $d\theta$.

- Remember Regularization.

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m (y_{(i)}, \hat{y}_{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_F^2$$

- Doesn't work with Dropout.

- Run at random initialization, perhaps again after some training.

$$\begin{cases} f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \Rightarrow O(\epsilon^2) \Rightarrow 0.01^2 \quad \checkmark \\ f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon) - f(\theta)}{\epsilon} \Rightarrow O(\epsilon) \Rightarrow 0.01 \end{cases}$$

Ch 2 Optimization Algorithms

1. Mini-batch Gradient Descent (Big Data)

Batch vs Mini-batch:

Vectorization allows to efficiently compute on m examples.

$$\vec{x} = [x^{(1)} \ x^{(2)} \dots x^{(1000)} \ | \ x^{(1001)} \dots x^{(2000)} \ | \ \dots \ | \ \dots x^{(m)}]$$

(n_x, m) $\underbrace{x^{(1)}}_{(n_x, 1000)} \quad \underbrace{x^{(2)}}_{(n_x, 1000)} \quad \underbrace{x^{(1000)}}_{(n_x, 1000)}$ $\underbrace{x^{(1001)}}_{(n_x, 1000)} \dots \underbrace{x^{(2000)}}_{(n_x, 1000)} \dots \underbrace{x^{(m)}}_{(n_x, 1000)}$

$$\vec{Y} = [y^{(1)} \ y^{(2)} \dots y^{(1000)} \ | \ y^{(1001)} \dots y^{(2000)} \ | \ \dots \ | \ \dots y^{(m)}]$$

(1, m) $\underbrace{y^{(1)}}_{(1, 1000)} \quad \underbrace{y^{(2)}}_{(1, 1000)} \quad \underbrace{y^{(1000)}}_{(1, 1000)}$ $\underbrace{y^{(1001)}}_{(1, 1000)} \dots \underbrace{y^{(2000)}}_{(1, 1000)} \dots \underbrace{y^{(m)}}_{(1, 1000)}$

If m=5 million \Rightarrow 5000 mini-batches of 1000 each.

Mini-Batch T: x^{T+1}, y^{T+1} .

Notation: $x^{(t)}, z^{(t)}, x^{(t+1)}, y^{(t+1)}$.

for t=1, ..., 5000

forward prop on $x^{(t)}$

$$z^{(t)} = W^{(t)} x^{(t)} + b^{(t)}$$

$$A^{(t)} = g^{(t)}(z^{(t)})$$

⋮

$$A^{(t)} = g^{(t)}(z^{(t)})$$

$$\text{Compute Cost } J^{(t)} = \frac{1}{1000} \sum_{i=1}^b L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_{l=1}^L \|W^{(l)}\|_F^2$$

Backward Prop to Compute Gradient ~~Cost~~ Cost $J^{(t)}$ ($x^{(t)}, y^{(t)}$)

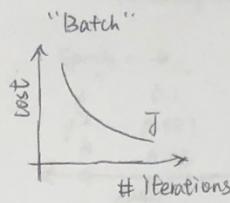
$$W^{(t)} = W^{(t)} - \alpha \cdot dW^{(t)} \quad b^{(t)} = b^{(t)} - \alpha \cdot db^{(t)}$$

不断训练 Training Set. 使它收敛在某个近似收敛值。

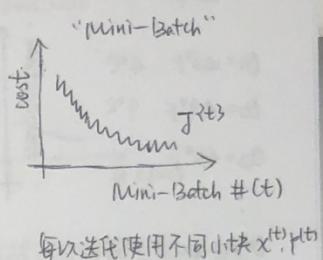
"Mini-Batch" faster than "Batch"

2. Understand Mini-batch Gradient Descent

1. Training with mini-batch Gradient Descent



每次迭代将遍历整个训练集



每次迭代使用不同小块 $x^{(t)}$ 趋势向下，会有许多噪声

2. Choose the Mini-batch Size

If Mini-batch size=m: Batch Gradient Descent.

$$(x^{(1)}, y^{(1)}) = (x, y)$$

If Mini-batch size=1: Stochastic Gradient Descent.

① Batch: Large Training Set \Rightarrow Too long per iteration.
Small Training Set is Okay. ($m \leq 2000$)

② Stochastic: Use a sample to renew gradient.

Lose speed-up from vectorization

每次只训练一个数据：Not efficient.

③ In-between: Fastest Learning.

- Vectorization.

- 不用等整个训练集都遍历完才运行梯度下降。每遍历一次训练集，可执行5000次梯度下降算法。

3. Typical Mini-batch Size

64	128	256	512	1024
2^6	2^7	2^8	2^9	2^{10}

Make mini-batch fit in CPU / GPU Memory: $x^{(t)}, y^{(t)}$

三. Exponentially Weighted Averages. 指数加权平均

Temperature in London:

$$\theta_1 = 40^{\circ}\text{F}, 4^{\circ}\text{C}$$

$$\theta_2 = 49^{\circ}\text{F}, 9^{\circ}\text{C}$$

$$\theta_3 = 45^{\circ}\text{F}$$

!

$$\theta_{180} = 60^{\circ}\text{F}, 15^{\circ}\text{C}$$

$$\theta_{181} = 56^{\circ}\text{F}$$

!

V_t 近似于 $\frac{1}{1-\beta}$ 天的均温

$\beta = 0.9 \approx 10$ days' temp.

$\beta = 0.98 \approx 50$ days' temp

$\beta = 0.5 \approx 2$ days' temp



$$V_0 = 0$$

$$V_1 = 0.9 V_0 + 0.1 \theta_1$$

$$V_2 = 0.9 V_1 + 0.1 \theta_2$$

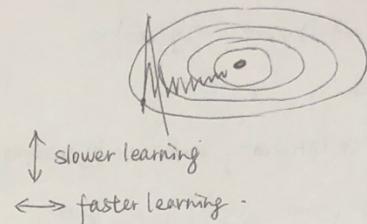
$$V_3 = 0.9 V_2 + 0.1 \theta_3$$

!

$$V_t = 0.9 V_{t-1} + 0.1 \theta_t$$

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

四. Gradient Descent with Momentum (动量)



上下振荡会减慢梯

度下降速度，无法使用

较大的 α .

(如果 α 太大，可能超调)

On iteration t : Compute dw, db on Mini-Batch.

$$Vdw = \beta \cdot Vdw + (1-\beta) dw$$

$$Vdb = \beta \cdot Vdb + (1-\beta) db$$

$$w = w - \alpha \cdot Vdw$$

↓ acceleration

$$b = b - \alpha \cdot Vdb$$

将动量项看作球沿碗
向下滚的速度。

β : friction.

V : velocity.



Hyperparameters: α, β . ($\beta = 0.9$ 是非常稳健的参数值).

五. RMSprop (均方根传递，“加速”)

↑ slow On iteration t : Compute dw, db
↔ fast on Mini-Batch

$$Sdw = \beta_1 Sdw + (1-\beta_1) dw^2 \leftarrow \text{small}$$

$$Sdb = \beta_2 Sdb + (1-\beta_2) db^2 \leftarrow \text{large}$$

$$w = w - \alpha \cdot \frac{dw}{\sqrt{Sdw + \epsilon}} \quad b = b - \alpha \cdot \frac{db}{\sqrt{Sdb + \epsilon}}$$

$\epsilon \approx 0 (\neq 0)$

Bias Correction:

$$V_0 = 0$$

$$V_1 = 0.98 V_0 + 0.02 \theta_1$$

$$V_2 = 0.98 V_1 + 0.02 \theta_2$$

$$= 0.98 \times 0.02 \theta_1 + 0.02 \theta_2$$

$$= 0.0196 \theta_1 + 0.02 \theta_2$$

$$\frac{V_t}{1-\beta^t}$$

$$t=2: |-\beta^2| = 1 - (0.98)^2 = 0.0396$$

$$\frac{V_2}{0.0396} = \frac{0.0196 \theta_1 + 0.02 \theta_2}{0.0396}$$

当 t 是非常大的时候，偏置校正不会影响操作。

当 t 是较小的时候，偏置校正给出一个更好的近似。

11. Adam Optimization Algorithm

(本质：put "Momentum" and "RMSprop" together)

已证明在许多不同种类的神经网络架构中都十分有效。

$$V_{dw} = 0, S_{dw} = 0, V_{db} = 0, S_{db} = 0$$

On iteration t : Compute dw, db on Mini-Batch.

$$\begin{aligned} V_{dw} &= \beta_1 \cdot V_{dw} + (1 - \beta_1) dw, \\ V_{db} &= \beta_1 \cdot V_{db} + (1 - \beta_1) db \end{aligned} \quad \left. \right\} \text{"Momentum" } (\beta_1)$$

$$\begin{aligned} S_{dw} &= \beta_2 \cdot S_{dw} + (1 - \beta_2) dw^2, \\ S_{db} &= \beta_2 \cdot S_{db} + (1 - \beta_2) db^2 \end{aligned} \quad \left. \right\} \text{"RMSprop" } (\beta_2)$$

$$V_{dw}^{corrected} = V_{dw} / (1 - \beta_1^t), V_{db}^{corrected} = V_{db} / (1 - \beta_1^t)$$

~~$S_{dw}^{corrected} = S_{dw} / (1 - \beta_2^t), S_{db}^{corrected} = S_{db} / (1 - \beta_2^t)$~~

$$w = w - \alpha \cdot \frac{V_{dw}^{corrected}}{\sqrt{S_{dw}^{corrected}} + \epsilon}, b = b - \alpha \cdot \frac{V_{db}^{corrected}}{\sqrt{S_{db}^{corrected}} + \epsilon}$$

Hyperparameters:

$$\left\{ \begin{array}{l} \alpha \text{ needs to be tuned. (尝试不同值比较效果)} \\ \beta_1 = 0.9 \text{ (动量算法, } dw \text{ 的移动平均/加权平均) \rightarrow \text{第一阶矩}} \\ \beta_2 = 0.999 \text{ (} dw^2, db^2 \text{ 的移动加权平均) \rightarrow \text{第二阶矩}} \\ \epsilon: 10^{-8} \text{ (Default). 都可以.} \end{array} \right.$$

设置一个固定的 α .

还要使它优化的好.

可真正加速算法运行

"Adam": Adaptive Moment Estimation (自适应矩估计)

12. Learning Rate Decay

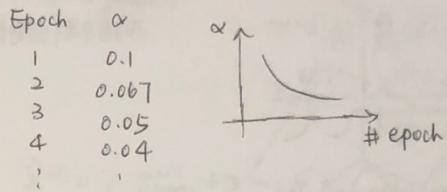
当用适量小样本迭代时，步长会逐渐向最小值靠近，但不会完全收敛至该点，(由于 α 取了固定值)

~~α~~ α Large. $\rightarrow \alpha$ small

(步长变小) \Rightarrow 最终围绕离极小值点更近的区域摆动

1 epoch = 1 pass through data.

$$\alpha = \frac{1}{(\text{decay rate} * \text{epoch.num})} \cdot \alpha_0$$

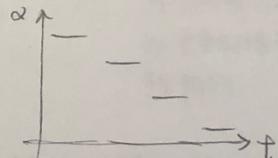


可尝试不同的超参数组合。
寻找一组效果好的数值 (α_0)

Other Learning Rate Decay Methods:

$$\alpha = 0.95^{\text{epoch.num}} \cdot \alpha_0 \quad (\text{exponentially decay})$$

$$\alpha = \frac{k}{\sqrt{\text{epoch.num}}} \cdot \alpha_0 \quad \text{or} \quad \frac{k}{\sqrt{t}} \cdot \alpha_0$$

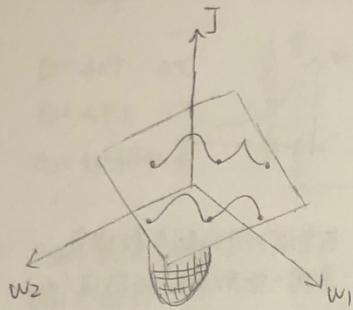


设置一个固定的 α .

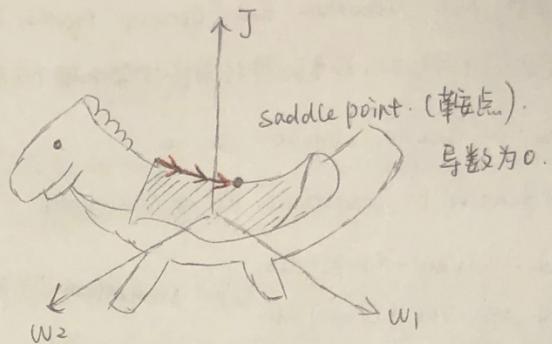
还要使它优化的好.

可真正加速算法运行

11. The problem of Local Optima.



容易陷入局部最优
而找不到全局最优.



Problem of plateaus:

“停滯帶區”. 导数长时间接近于0的一段区域.

Unlikely to get stuck in a bad local optima.
plateaus can make learning slow

“Adam Algorithm”: 加快停滯區向下移动.

Ch.3 Tuning Process

- Hyperparameter Tuning

1. Tuning Process

对于超参数，Try random values. 而不是在网格中规则抽样。
将更充分地为最重要的超参数尝试多种可能组合。

⇒ 区域限定 在某区进行更密集的抽样 (区域搜索)

2. Using an appropriate scale to pick hyperparameters.

随机抽样

在 $a \sim b$ 范围内均匀随机取样

3. Hyperparameters tuning in Practice: Pandas vs. Caviar

① Re-test hyperparameters occasionally: 确保这些数值依旧是最优解

② Babysitting one model

③ Training many models in parallel

Give some intermediate values in NN:

$$\underbrace{z^{(1)} \dots z^{(m)}}_{z^{(i)}} \quad \mu = \frac{1}{m} \sum_{i=1}^m z^{(i)}$$

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (z^{(i)} - \mu)^2$$

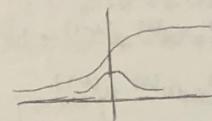
$$z_{\text{norm}} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\hat{z}^{(i)} = \gamma \cdot z_{\text{norm}} + \beta \quad (\gamma, \beta: \text{can be learned from model})$$

$$\text{If } \gamma = \sqrt{\sigma^2 + \epsilon}$$

$$\beta = \mu$$

$$\text{Then } \hat{z}^{(i)} \approx z^{(i)}$$



隐藏层归一化后并不一定是均值为0. 方差为1.

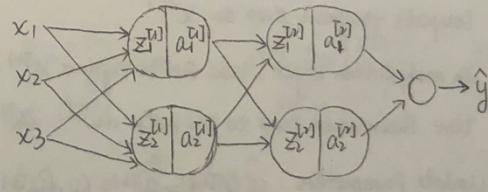
如果激活函数是 sigmoid.

可能希望它们有更大方差.

以便更好地利用S函数非线性特性.

2. Fitting Batch Norm into a Neural Network

① Adding Batch Norm to a Network:



$$x \xrightarrow{w^T, b^T} z^T \xrightarrow{\mu^T, \sigma^T} \hat{z}^T \xrightarrow{\text{Batch Norm (BN)}} a^T = g^T(\hat{z}^T)$$

$$w^T, b^T \xrightarrow{z^T} \hat{z}^T \xrightarrow{\mu^T, \sigma^T} \hat{z}^T \xrightarrow{\text{BN}} a^T \dots$$

Parameters:

$$\left. \begin{array}{l} w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}, \dots, w^{(L)}, b^{(L)} \\ \beta^{(1)}, \gamma^{(1)}, \beta^{(2)}, \gamma^{(2)}, \dots, \beta^{(L)}, \gamma^{(L)} \end{array} \right\} d\beta^{(L)} \quad \text{BN} = \beta^{(L)} - \alpha \cdot d\beta^{(L)}$$

tf.nn.batch_normalization.

二. Batch Normalization

1. Normalize Activations in a Network.

Normalize Inputs to speed up Learning:

$$x_1, w, b \xrightarrow{O} \hat{y}$$

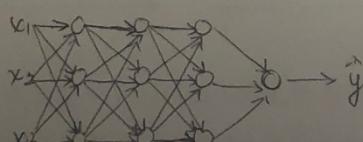
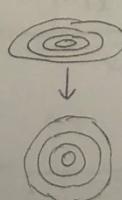
$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$x = x - \mu$$

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)})^2$$

$$x = x / \sigma^2$$

(Logistic Regression.)



② Working with Mini-Batches:

$$x^{(i)} \xrightarrow{w^{(i)}, b^{(i)}} z^{(i)} \xrightarrow{\beta^{(i)}, \gamma^{(i)}} \hat{z}^{(i)} \xrightarrow{g^{(i)}(\hat{z}^{(i)})} a^{(i)}$$

$$\underbrace{w^{(i)}, b^{(i)}}_{\text{parameters}} z^{(i)} \rightarrow \dots$$

$$x^{(i)} \xrightarrow{\text{BN}} z^{(i)} \xrightarrow{\beta^{(i)}, \gamma^{(i)}} \hat{z}^{(i)} \rightarrow \dots$$

...

parameters: $w^{(i)}, b^{(i)}, \beta^{(i)}, \gamma^{(i)}$ (维度: $n^{(i), 1}$)

$$z^{(i)} = w^{(i)} a^{(i-1)} + b^{(i)}$$

$$z^{(i)} = w^{(i)} a^{(i-1)}$$

$$\hat{z}^{(i)}$$

$$\hat{z}^{(i)} = \gamma^{(i)} \cdot z_{\text{norm}}^{(i)} + \boxed{\beta^{(i)}}$$

BN算法使层级中各个 $\hat{z}^{(i)}$ 的均值为0. 因此无需保留 $b^{(i)}$ 而是用 $\beta^{(i)}$ 来代替. 控制最终偏移量影响的参数.

③ Implement Gradient Descent:

for $t=1 \dots$ Mini-Batches

Compute forward prop on $x^{(t)}$

In each hidden layer, use BN to replace $z^{(i)}$ with $\hat{z}^{(i)}$

Use Backward prop to compute $d\hat{z}^{(i)}, d\hat{z}^{(i)}, d\beta^{(i)}, d\gamma^{(i)}$

Update parameters: $w^{(i)} = w^{(i)} - \alpha \cdot d\hat{z}^{(i)}$

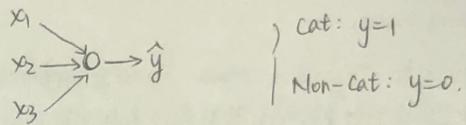
$\beta^{(i)} = \beta^{(i)} - \alpha \cdot d\beta^{(i)}$

$\gamma^{(i)} \dots$

同样适用于 Momentum · RMSprop · Adam 的梯度下降算法.

3. Why does Batch Norm work?

Learning on shifting Input Distribution:



BN: 减少了隐藏单元值分布的不稳定性.

(减少输入值变化所产生的问题)

神经网络的后层, 可以有更加稳固的基础)

⇒ 后面层适应前面层变化的力量被减弱.

Batch Norm as Regularization:

- Each mini-batch is scaled by the mean/ variance computed on just that mini-batch.
- This adds some noise to the values $\hat{z}^{(i)}$ within that mini-batch. So similar to dropout, it adds some noise to each hidden layer's activations.
- This has a slight regularization effect.

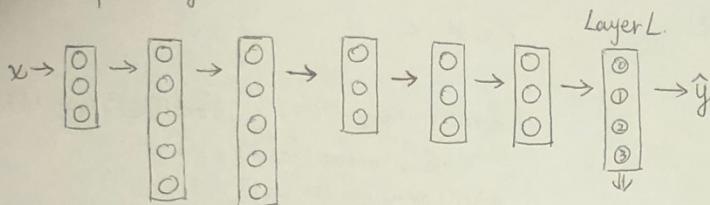
4. Batch Norm at test time.

$$\mu = \frac{1}{m} \sum_i z^{(i)}, \quad \sigma^2 = \frac{1}{m} \sum_i (z^{(i)} - \mu)^2$$

$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}, \quad \hat{z}^{(i)} = \gamma \cdot z_{\text{norm}}^{(i)} + \beta.$$

3. Multi-Class Classification

1. SoftMax Regression.



$$z^{(l)} = w^{(l)} a^{(l-1)} + b^{(l)}$$

Activation Function: $t = e^{z^{(l)}}$

$$a^{(l)} = \frac{e^{z^{(l)}}}{\sum_{j=1}^C t_j}, \quad a_j^{(l)} = \frac{t_j}{\sum_{j=1}^C t_j}$$

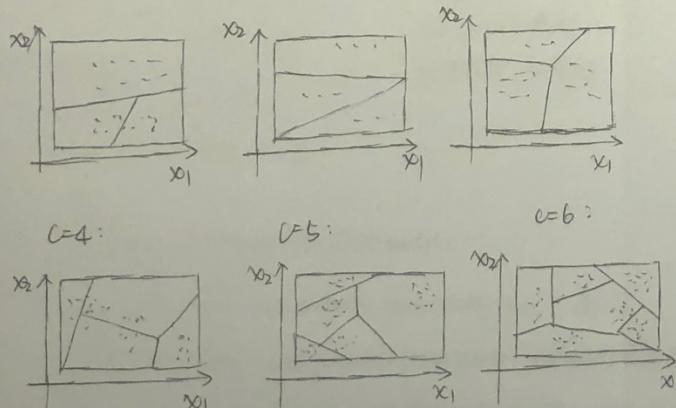
$$a^{(l)} = g^{(l)}(z^{(l)})$$

$$z^{(l)} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix} \cdot \frac{4}{\sum_{j=1}^4 t_j} = 176.3$$

$$a^{(l)} = \frac{t}{176.3}$$

2. Softmax Example.

$C=3$: (决策边界是线性的，不同的代价函数区分不同类别)



这些图展示了 SoftMax 分类器在没有隐藏层时，可以做什么。有很深的网格，可用更复杂的非线性决策平面来区分不同类别。

3. Training a Softmax Classifier.

① Understand Softmax:

$$z^{(l)} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix}$$

$$\begin{aligned} a^{(l)} &= g^{(l)}(z^{(l)}) = \begin{bmatrix} e^5/(e^5+e^2+e^{-1}+e^3) \\ e^2/(e^5+e^2+e^{-1}+e^3) \\ e^{-1}/(e^5+e^2+e^{-1}+e^3) \\ e^3/(e^5+e^2+e^{-1}+e^3) \end{bmatrix} \\ &= \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix} \quad \text{"Hard Max"} \quad \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \end{aligned}$$

Softmax Regression generalizes logistic regression to C classes.

If $C=2$, softmax reduces to logistic regression.

$$a^{(l)} = \begin{bmatrix} 0.842 \\ 0.158 \end{bmatrix}$$

② Loss Function:

$$\hat{y} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \rightarrow \text{cat.} \quad a^{(l)} = \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix} \quad C=4 \quad (4,1)$$

$$L(\hat{y}, y) = -\sum_{j=1}^C y_j \log \hat{y}_j$$

small

$$-y_2 \log \hat{y}_2 = -\log \hat{y}_2. \text{ Make } \hat{y}_2 \text{ big.}$$

$$J(w^{(l)}, b^{(l)}, \dots) = \frac{1}{m} \cdot \sum_{i=1}^m L(\hat{y}_i, y_i)$$

$$\begin{aligned} Y &= [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}] \quad \hat{Y} = [\hat{y}^{(1)} \ \hat{y}^{(2)} \ \dots \ \hat{y}^{(m)}] \\ &= \begin{bmatrix} 0 & 0 & 0 & \dots \\ 1 & 0 & 0 & \dots \\ 0 & 1 & 0 & \dots \\ 0 & 0 & 1 & \dots \end{bmatrix} \quad = \begin{bmatrix} 0.3 & \dots \\ 0.2 & \dots \\ 0.1 & \dots \\ 0.4 & \dots \end{bmatrix} \\ &\quad (4, m) \quad (4, m) \end{aligned}$$

③ Gradient Descent with SoftMax.

Forward prop: $\underline{z}^m \rightarrow a^m = \hat{y} \rightarrow L(\hat{y}, y)$

Backward prop: $d\underline{z}^m = \hat{y} - y \quad \frac{\partial J}{\partial \underline{z}^m}$

IV. Programming Frameworks.

1. Deep Learning Frameworks

Caffe / Caffe 2
CNTK
DL4J
Keras
Lasagne
mxnet
Paddle Paddle
TensorFlow
Theano
Torch

★ Choose DL Frameworks:
 ① Ease of programming
 (development & deployment)
 ② Running Speed.
 ③ Truly Open.
 (Open Source with good governance)

session.run(train). 在 train 上运行下一步梯度下降
 print(session.run(w))

Output: 0.1

for i in range(1000): 将迭代运行 1000 次梯度下降
 session.run(train)
 print(session.run(w))

Output: 4.99999. } Similar.
 w(optimal) = 5.

矩阵:

```
coefficients = np.array([[1.], [-10.], [25.]])  

w = tf.Variable(0, dtype=tf.float32)  

x = tf.placeholder(tf.float32, [3, 1])  

cost = x[0]*x[0]*w**2 + x[1]*x[1]*w + x[2]*w  

session.run(train, feed_dict={x: coefficients})
```

Output: same with ~~w~~. Before

2. TensorFlow

Motivate Problem:

$$J(w) = w^2 - 10w + 25 = (w-5)^2$$

(cost function) $w=5$.

```
import numpy as np
```

```
import tensorflow as tf
```

$w = tf.Variable(0, dtype=tf.float32)$ 定义参数. (尝试优化的值)

$cost = tf.add(tf.add(w*w), tf.multiply(-10, w)), 25)$ 定义代价函数 J.

$train = tf.train.GradientDescentOptimizer(0.01).minimize(cost)$ 最小化代价函数. \Leftarrow 只需更改此行.

$init = tf.global_variables_initializer()$

$session = tf.Session()$ 启动一个 TensorFlow 的 Session.

$session.run(init)$ 初始化全局变量.

$sess.run(w)$. 对变量 w 求值. ($print(session.run(w))$) \Rightarrow Output: 0.0

等同于 with tf.Session() as session:
 session.run(init)
 print(session.run(w))

Name: Xiaofan Wang

Structuring Machine Learning Projects

Ch 1

1. Introduction to ML Strategy.

1. Why ML Strategy.

学习更高效的构建机器学习系统的方法

Example: 开发猫分类模型 \Rightarrow 系统准确率 90%

Aim: 提升准确率

Ideas:

- Collect more data
- Collect more diverse training set
- Train Algorithm longer with gradient descent.
- Try Adam instead of gradient descent
- Try bigger or smaller network.
- Try dropout.
- Add L2 regularization.

Network Architecture:

Activation Functions	{	# Hidden Units
-		

★ 分辨哪些方法值得尝试，哪些应舍弃（节省时间）。

2. Orthogonalization 正交化

One of the challenges with building ML systems:

So many things you could try / change. (有许多可尝试和改变的)

e.g. so many hyperparameters you could tune. 许多可调的超参数。

最高效的机器学习人员：非常清楚需要调整什么

\Rightarrow Orthogonalization. 来实现一个预期效果。

$$F1\text{-Score} = \frac{2}{P+R} \quad (\text{调和平均值})$$

① Classifier Precision Recall F1-Score 权衡 PR

Classifier	Precision	Recall	F1-Score	权衡 PR
A	95%	90%	92.4%	✓
B	98%	85%	91%	

一个好的验证集和单一量化评估指标

可以加速迭代，让 ML-Algorithm 不断改进。

② Another Example. (error rate)

Algorithm	US	China	India	Average
A	3%	7%	5%	6%
B	5%	6%	5%	6.5%
C	2%	3%	4%	(3.5%) C Better.
D	5%	8%	7%	5.25%
E	4%	5%	2%	3.75%
F	7%	11%	8%	9.5%

C Better.

选择 C 算法进行迭代



调整图像的旋钮

$\downarrow \times 0.1$ (每个旋钮都有一个实际功能)

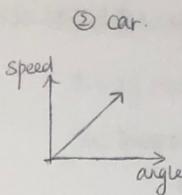
$\leftarrow \times 0.3$ 假设一个旋钮可同时控制

$\leftarrow \square \rightarrow \times 1.7$ $a + 0.3b - 1.7c - \dots$

$\square \times 0.8$ 那用户几乎无法调整电视菜单居中

$\square \dots$ \Rightarrow Orthogonalization: 确保每个旋钮只可进行一个参数的调整

使画面居中



$$0.3 \times \text{angle} - 0.8 \times \text{speed}$$

$$2 \times \text{angle} + 0.9 \times \text{speed}$$

正交控制：可与真正想控制的事物保持一致。

③ Chain of assumptions in ML. (正交性原理)

调整系壳旋钮，保证 4 点：

- Fit training set well on cost function (performance).
- Fit Dev Set well on cost function.
- Fit test set well on cost function
- Perform well in real world.

2. Set up your goal

1. Single Number Evaluation Metric.

(设置一个单一的量化评估指标)

①

Classifier	Precision	Recall	F1-Score	权衡 PR
A	95%	90%	92.4%	✓
B	98%	85%	91%	

一个好的验证集和单一量化评估指标

可以加速迭代，让 ML-Algorithm 不断改进。

② Another Example. (error rate)

Algorithm	US	China	India	Average
A	3%	7%	5%	6%
B	5%	6%	5%	6.5%
C	2%	3%	4%	(3.5%) C Better.
D	5%	8%	7%	5.25%
E	4%	5%	2%	3.75%
F	7%	11%	8%	9.5%

C Better.

选择 C 算法进行迭代

2. Satisficing and Optimizing Metrics. (满足指标, 优化指标)

① Another Cat Classification Example

Classifier	Accuracy	Running Time.
A	90%	80ms
B	92%	95ms
C	95%	1500ms.
	↓	↓
optimizing		Satisfying

$$cost = \text{Accuracy} - 0.5 \times \text{Running Time.}$$

Maximize Accuracy.

satisfy: $\text{Running Time} \leq 100\text{ms}$ (满足条件即可)

对准确率和运行时间, 作权衡或通盘考虑.

N metric: [1 optimizing.

[N-1 Satisfying (达到阈值即可)

② Wake Words / Trigger Words.

- e.g. | Amazon Echo: Alexa
| Google: OK Google. \Rightarrow 用于告诉某个声控设备你想表达的意思.
| Apple: Hey Siri
| 度: 你好度.
 • 触发词监测系统的准确度 (即实际唤醒系统的能力).
 • False Positive: 随机唤醒的可能性.

3. Train / Dev / Test Distributions

Cat Classification Dev / Test Sets.

Regions:

US	Dev.	结果不准确, 相当于换了BIN.
UK		数据来自不同的分布.
Other Europe	Test.	应使8个地区数据随机打乱并分为 Dev / Test Sets.
South America		(这样 Dev / Test Dev 来自相同分布)
India	Test.	
China		
Other Asia	Test.	
Australia		

True Story: (different distribution).

中产 [Optimizing on dev set on loan approvals for medium income zip codes.
低收入 [Tested on low income zip codes.

(不准确)

Guideline: 来自相同分布

choose a dev set and test set to reflect data you expect to get in the future and consider important to do well on.

(A or B.)

4. Size of Dev / Test Sets. [Dev Set: 训练不同想法
Test Set: 对最终分类器进行评估

① Old way of splitting Data: (不再适用)

70%	30%	
Train	Test.	
60%	20%	20%
Train	Dev	Test.

样本量不大时, 很适用

100
1000
10000

Modern:

98%	1%	1%
Train	Dev	Test.

数据量很大时, (100000)

没必要划分 7/3.
6/2/2.

(确保 Dev. Test Set 的大小足够达到目的)

因为深度学习算法极度需要数据, 所以训练集占很大比例.

② Size of Test set

Set test set to be big enough to give high confidence in the overall performance of your system. (足够能保证对系统整体性能评估的高置信度即可)

除非对于“最终系统性能”有非常精确的测量, 否则测试集不需要上百万个样例.

③ * Only Train / Dev Set.

在一些应用中, 不需对最终系统的整体性能评估有很高的置信度. \Rightarrow 只需 Train, Dev Sets. No Test Set.

Train	Dev
-------	-----

如果 Dev Set 很大, 过拟合得厉害.
可取此法, 但不建议.

5. When to change dev/test sets and metrics.

① Cat Dataset ~~评估~~

Metric: classification error

Algorithm A: 3% error. (pornographic...)
Algorithm B: 5% error.

⇒ Metric + Dev Set: prefer A.

You / Users: prefer B.

⇒ 修改评估指标. 修改 Dev/Test Set.

只是定义评估指标的一中方法，统计分类错误的指标： $\frac{1}{M_{dev}} \cdot \sum_{i=1}^{M_{dev}} L(\hat{y}^{(i)}, y^{(i)})$
 考虑到 pornographic 权重 $w^{(i)}$: $\frac{1}{\sum_i w^{(i)}} \cdot \sum_{i=1}^{M_{dev}} w^{(i)} L(\hat{y}^{(i)}, y^{(i)})$.
 $w^{(i)} = \begin{cases} 1 & \text{if } x^{(i)} \text{ is non-pornographic} \\ 10 & \text{if } x^{(i)} \text{ is pornographic} \end{cases}$

② Orthogonalization for cat pictures: anti-porn.

- So far we've only discussed how to define a metric (Place Target) to evaluate classifiers. 确定一个指标以衡量分类器在目标上的性能。
- Worry Separately about how to do well on this metric. 单独考虑如何在这个指标上得到很好地性能. (Aim / Shoot at the Target).

$$J = \frac{1}{m} \cdot \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

$$J = \frac{1}{\sum_i w^{(i)}} \cdot \sum_{i=1}^{w^{(i)}} w^{(i)} L(\hat{y}^{(i)}, y^{(i)})$$

③ If doing well on your metric + dev/test set does not correspond to doing well on your application, change your metric and dev/test set.

对于团队工作: 迅速确定评估指标. Dev/Test Set

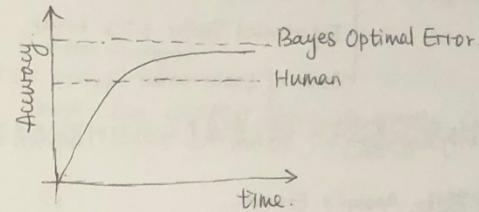
可以进行迭代.

如果效果不佳, 之后可再做修改.

三. Comparing to Human-Level Performance.

1. Why Human-Level Performance.

ML 超越 Human: ① 算法效果更好
② 解决问题的效率高.



- ① ML 超越 Human-level 后, 可得到更好的效果, 但增长速度变慢
- ② 可无限逼近最佳理论水平 (Bayes Optimal Error). 最佳理论误差值

③ Humans are quite good at a lot of tasks. So long as ML is worse than humans, you can:

- Get labeled data from humans.
- Gain insight from manual error analysis: Why did a person get this right?
- Better analysis of bias/variance

2. Avoidable Bias.

	Bias		
Humans	1%	(B)	Avoidable Bias
Training error	8%	7% (B)	0.5% (图片太糊)
Dev error	10%	2% (V)	8% 2% (UV)

Focus on "Bias"

↓
Train bigger Neural Network.

Train longer time

↓
Improve performance of model on Training set.

Focus on "Variance"

改进模型: 误差小

到与 Bayes 误差齐平

小于 Bayes 误差不可能
除非过拟合.

3. Understand Human-Level Performance.

① Human-Level Error as a proxy for Bayes Error.

Medical Image Classification Example:

Suppose:

$$\left\{ \begin{array}{l} \text{Typical Human } 3\% \text{ error} \\ \text{Typical Doctor } 1\% \text{ error} \\ \text{Experienced Doctor } 0.7\% \text{ error} \end{array} \right.$$

Team of Experienced Doctors 0.5% error.

如何定义人类误差：将人类误差当成贝叶斯误差的简化值或估计值。

② Error Analysis Example.

Human (proxy for Bayes Error)	$\left\{ \begin{array}{l} 1\% \\ 0.7\% \\ 0.5\% \end{array} \right.$	$\left\{ \begin{array}{l} 1\% \\ 0.7\% \\ 0.5\% \end{array} \right.$	$(0.7\% \text{ 更难改进})$
\downarrow Avoidable Bias	$\downarrow 4\% - 4.5\%$	$\downarrow 0.05\%$	$\downarrow 0.2\%$
Training Error	5%	1%	0.7%
\downarrow Variance	$\downarrow 1\%$	$\downarrow 4\%$	$\downarrow 0.1\%$
Dev Error	6%	5%	0.8%
\downarrow	\downarrow	\downarrow	\downarrow
Bias		Variance	<u>偏差 方差都可改进</u>

4. Surpass Human-Level performance.

- Online Advertising.
- Product Recommendation.
- Logistics (predicting transit time)
- Loan Approvals.

单个监督学习 > 人类水平

Structured Data (Lots of Data) = ML > Human-Level

Not Natural Language Processing. (人类更擅长 NLP)

NLP 一些突破性例外：

$\left\{ \begin{array}{l} \text{Speed recognition} \\ \text{Image recognition} \end{array} \right.$	$\text{Medical: EGG, Skin cancer}$
---	------------------------------------

5. Improve Model Performance.

① Two fundamental assumptions of supervised learning:

$\left\{ \begin{array}{l} \text{You can fit the training set pretty well.} \\ \text{The training set performance generalizes pretty well to the dev/test set.} \end{array} \right.$

② Reduce Avoidable Bias and Variance.

Human-Level	\uparrow	Train bigger model				
Avoidable Bias	\downarrow	Train longer/better optimization algorithms.				
Training Error	\Rightarrow	NN architecture/hyperparameter search.				
Variance	\downarrow					
Dev Error	\Rightarrow	<table border="0"> <tr> <td>More Data</td> </tr> <tr> <td>Regularization: L2, dropout.</td> </tr> <tr> <td>Data Augmentation</td> </tr> <tr> <td>NN architecture/hyperparameters search.</td> </tr> </table>	More Data	Regularization: L2, dropout.	Data Augmentation	NN architecture/hyperparameters search.
More Data						
Regularization: L2, dropout.						
Data Augmentation						
NN architecture/hyperparameters search.						

Chapter 2.

1. Error Analysis

1. Carrying out error analysis.

① Look at dev examples to evaluate ideas.

Should you try to make your cat classifier do better on dogs?

Error Analysis: { Get ~100 mislabeled dev set examples
(手动检测) Count up how many are dogs.

5 / 100. 5%. (error: 10% → 9.5%)



并不是最好的利用时间的方式。

50 / 100. 50% (error: 10% → 5%)



值得付出努力：专注于减少被错误标识的狗。

② Evaluate Multiple Ideas in Parallel.

Ideas for Cat Detection: (多个改进猫检测器的思路).

- Fix pictures of dogs being recognized as cats.
- Fix great cats (non-potters...) being misrecognized.
- Improve performance on blurry images.

Image	Dog.	Great Cats	Blurry	Instagram	Comments
1	✓			✓	Pitbull
2			✓	✓	
3		✓	✓		Rainy Day at Zoo.
⋮	⋮	⋮	⋮	⋮	
% of total (分类错误 的百分比)	8%	43%	61%	12%	提高准确率的幅度大。

结果：不同类别的错误，有多少处理价值。

总结：错误分析 → 找到一套在验证集中错误标识的样本。

→ False positive. False Negative

计算不同类别的误判个数 (判断优先级)

→ 可能有新的错误类别。新策略。

2. Cleaning up Incorrectly labeled data (Output y)

① Training Set.

• DL Algorithms are quite robust to random errors

In the training set, especially in lots of data.

• But Not robust to systematic errors (习惯性错误)

② Error Analysis.

Image	Dog	Great Cat	Blurry	Incorrectly Labeled	Comments
-------	-----	-----------	--------	---------------------	----------

98				✓	taker missed cat in background.
99				✓	Drawing of a cat.
100				✓	Not a real cat.
% of total	8%	43%	61%	6%	

Overall Dev Set Error = 10%

Errors due to incorrect labels: 0.6% (X)	0.6% (X)
--	----------

Errors due to other causes: 9.4%	1.4%
----------------------------------	------

(纠正无太意义) (纠正有效)

Goal of Dev Set: To help you select between two classifiers A & B.

③ Correcting Incorrect Dev/Test Set Examples.

• Apply same process to your dev and test sets to make sure they continue to come from the same distribution.

(Train and dev/test set may come from slightly different distributions.)

• Consider examining examples your algorithms got right as well as ones it got wrong.
(98%) (2%)

当一些标签错误的样本被纠正时，很可能预测会变得不准。

总结：当带领团队工作时，希望理解这个机

器犯的错。→ 手动检查错误。

⇒ 可更好地计划 Next Step.

3. Build your first system quickly, then iterate.

① Speech Recognition Example. 语音识别系统.

- Noisy Background = cafe noise, car noise
- Accented Speech
- Far from microphone.
- Tiny children's speech
- Stuttering.

可以有很多方面来改善语音识别系统。

②

- Set up dev/test set and metric 建立开发集/测试集和评价指标 (target)
- Build initial system quickly.
找到训练集，训练它然后看结果。
开始观察并理解你的系统，对开发/训练集如何
价值和度量指标是多少。
- Use Bias/Variance Analysis & Error Analysis
to prioritize next steps.

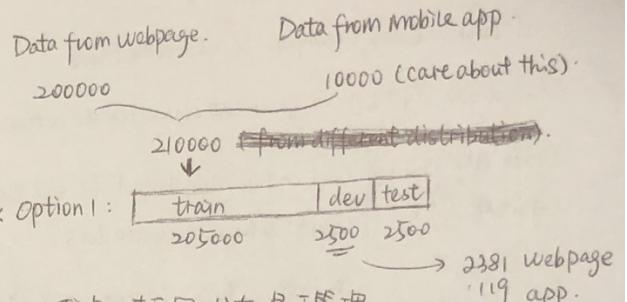
总结：快速建立你的第一个系统，然后迭代。
(不适用于你已经经验丰富的领域)
(不适用于已有很多可参考的学术论文)
↓
比如：人脸识别
所依托论文：从一开始建立一个复杂系统。

但如果是研究一个新问题，不要考虑太多。
↓
不要将系统做的太复杂
快速建立一个系统，然后使用它。
来帮助你确定一个优先级来改善你的系统。

2. Mismatched training and dev/test data.

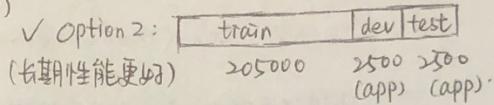
1. Training and testing on different distributions.

① cat app example.



优点：来自同一分布，易于管理。

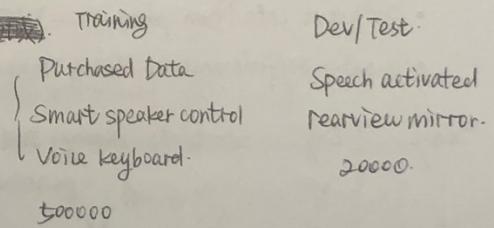
缺点：Dev Set 中分布比例失衡，不能重点关注 App。



优点：瞄准目标。

缺点：train set 与 dev/test set 分布不同。

② Speech recognition example.



Train	Dev	Test
500K.	10K	10K

Train	Dev	Test
510K.	5K	5K

2. Bias and Variance with Mismatched Data Distributions

① Cat classifier example: (Assume humans get $\approx 0\%$ error).

Training Error: 1%

Dev Error: 10%

\Rightarrow 如果 Train, Dev Set 来自同一分布, 说明方差过大.
虽然 Train Set 训练的很好, 但未能较好泛化到 Dev Set.

\Rightarrow 如果 Train, Dev Set 来自不同分布.
只因高分辨率的图片, 造成训练集难度较低.
但 Dev Set 较难, 不是方差问题, 而是 Dev Set 难以精准区分.

② Training-dev Set: Same distribution as training set.
but not used for training.

0/1/1/1/1 Train / / / / / T-D dev Test

NN	(Human)		
Training Error 1%	0%	Avoidable Bias.	0% { Avoidable Bias.
Train-Dev Error 9%	1%	10%	10% { Avoidable Bias.
Dev Error 10%	1.5%	11%	11% { Data mismatch.
		12%	20% { Data mismatch.
High Variance.	Data. (关联问题)	Avoidable Bias.	Bias + Data mismatch.
	Mismatch.	Mismatch.	Mismatch.

③ Bias / Variance on Mismatched training and Dev/Test Sets.

Human Level	4%	Avoidable Bias	4%
Training Set Error.	7%	Variance	7%
Training-Dev Set Error	10%	Mismatch.	10%
Dev Error.	12%	Data mismatch.	6%
Test Error.	12%	Degree of overfitting to Dev Set.	6%.

如果 Dev, Test Set 性能差太多, 将神经网络调得
得太偏向 (Overtuned) 于开发集了.

④ More general formulation:

Human Level	General (语音识别) rearviewmirror
Error on Examples trained	4% 6% { Bias
Error on not examples trained	Train Set 7% 6% { Var
Train-Dev Set	10% { Dev/Test
	6% { Data mismatch.

总结: 解决 Data mismatch.

Data Mismatch.

实际上并没有很好的(系统地)处理方法

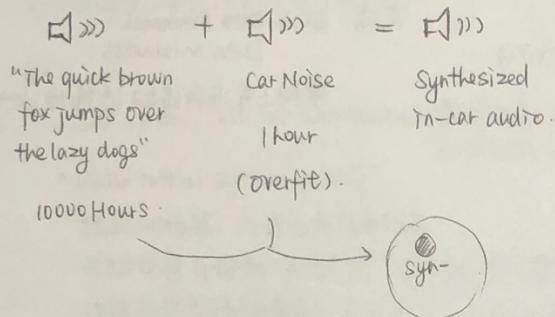
特别: 语音识别.

从上到下
应该是↑的.

3. Address Data Mismatch.

- ① Carry out manual error analysis to try to understand difference between training and dev/test sets.
- ② Make training data more similar.
Or collect more data similar to Dev/Test sets.

⇒ ④ Artificial Data Synthesis 人工合成 (应用: 语音识别).



总结: 如果遇到 ③ Data Mismatch. ⇒ Error Analysis.

比较 Train Set. vs. Dev Set

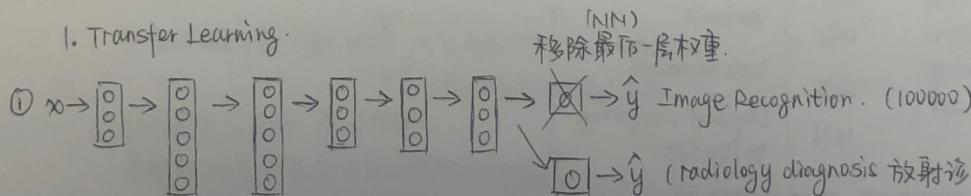
了解两种数据分布的不同。
可设法获得更多像开发集的训练数据。

。 人工合成. 可应用于语音识别.

使用时注意. 是否意外地模拟了所有可能样本的子集的数据.

三. Learning from multiple tasks.

1. Transfer Learning.



(x, y) . 放射影像.

⇒ 然后重新在新的数据集上

训练该 NN (新放射数据) (数据多)

⇒ 也可只重新训练最后一层权重 $w^{(l)}$
~~而保留前几层~~. (数据少)

移除最后一层权重.

最后一层为 NN 最后一层创建

一个新的随机初始化的权重. ⇒ 使用新建的输出层

来进行放射结果诊断.

当我们在做 Image Recognition 时. 训练了所有常用的神经网络参数/权值
从而做出图像识别的预测.

2. Multi-task Learning (每个任务将有助于其它任务).

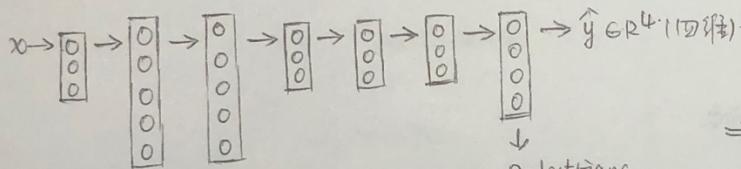
① Simplified autonomous driving example:

Image	Pedestrians	$y^{(i)}$
$x^{(i)}$	Cars	0
	Stop Signs	1
	Traffic Lights	1
		0

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

(4, m).

② Neural Network Architecture:



为了训练这个神经网络.

需要定义~~是~~损失函数:

$$\text{Loss: } \hat{y}^{(i)} \rightarrow \frac{1}{m} \sum_{j=1}^m L(\hat{y}_j^{(i)}, y_j^{(i)})$$

↓
usual Logistic Loss

$$y_j^{(i)} \log \hat{y}_j^{(i)} + (1 - y_j^{(i)}) \log (1 - \hat{y}_j^{(i)})$$

多
标
签
↓
Pedestrians
Cars
Stop Signs
Traffic Lights

⇒ 与 Softmax 回归不同.

(一个标签).

⇒ 该例是“多标签”. 一张图片解决4个问题.

试图告诉你每个图像中的4个对象.

或训练4个独立的NN.

但不同的输出之间, 前面的特征可共享.

⇒ $Y = \begin{bmatrix} 1 & 0 & ? & ? \\ 0 & 0 & ? & ? \\ 1 & 1 & 0 & ? \\ 1 & 1 & 0 & ? \end{bmatrix}$ 存在许多“?”的 Dataset.

你需要训练你的算法.

③ When Multi-Task Learning Makes Sense.

⇒ Training on a set of tasks that could benefit from having shared lower-level features. (可共享低层特征).

⇒ Usually: Amount of Data you have for each task is quite similar. (每个单任务数据量相似).

{ Transfer Learning: A (1000000) → B (1000)

Multi-Task Learning:
 ↓
 A1 1000
 A2 1000
 :
 An 1000

如果某个任务1000, 其它任务合计有远超1000的数据, 那么, 其它任务可帮助你在最后的任务表现更好.

⇒ Can train a big enough NN to do well on all tasks.

很大的NN使所有任务都有高准确度.

为每个任务单独训练一个NN.

(如果NN不够大, 会损害准确度.)

总结: 多任务学习应用的比迁移学习少得多.

多任务学习可训练一个NN来完成多任务. 这可使得性能比单独执行任务要好.

{ 多任务学习

迁移学习: 数据少, 因此应用多.

(例外: 计算机视觉物体检测).

实际上迁移学习应用更多.

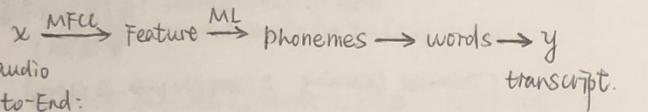
④ End-to-End Deep Learning. (端到端深度学习)

1. What is end-to-end deep learning.

一些数据处理系统，或者是由多个阶段组成的学习系统。

通常可将其替换为单个NN

① Pipeline:



② End-to-End:



End-to-End Deep Learning: 在最新的训练数据中。

(需要大量数据)

直接学习 $x \rightarrow y$ 的映射。

(才有好结果)

当数据少时，传统的 pipeline 同样有效。

③ audio - - -> phonemes - - -> transcript.

不彻底的端到端的学习。

④ Face Recognition.

Image (x)

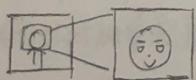


→ Identity (y)

函数映射。(一个模型同时学习所有步骤)。

但当人位置改变时不再适用。

只能获取少许 (x, y)。



→ Identity

将图片放入NN。

放大该人的脸部。

进行识别。

可得到许多 (x, y) 样本。

⑤ More examples.

• Machine Translation: English → Text Analysis → ... → French.

(x, y) 可得大量数据。 English → - - - → French.

English ⇔ French.

• Estimating Child's Age: Image → Bones → Age.

Image → - - - → Age. (不行)



2. Whether to use end-to-end learning.

① Pros and Cons of end-to-end learning.

Pros: { Let the data speak.

{ less Hand-designing of components needed

Cons: { May need large amount of data.

{ Exclude potentially useful hand-designed components

② Applying End-to-End Deep Learning.

Key Question: Do you have sufficient data to learn
a function of the complexity needed
to map x to y .

CNN Name: xiaofan Wang

Ch 1 Foundations of Convolutional Neural Network.

1. Computer Vision

1. Image Classification: Cat (0/1)
- Object Detection
- Neural Style Transfer: Figure mixed with picture.

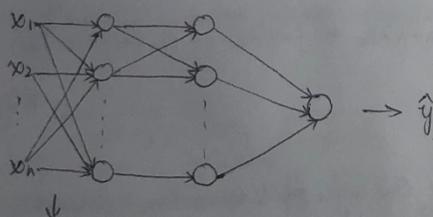
2. One challenge in CV: inputs can get really big.

e.g. picture: $64 \times 64 \times 3$ (3 color channels) $\underline{\text{RGB}}$

the input features has dimension 12288

e.g. Larger Image: 1000 pixel \times 1000 pixel \times 3

\rightarrow 3 million input features \rightarrow 3 million dimensional



1st Hidden layer: 1000 hidden units.

Total weights: Matrix \vec{W}_1 (1000, 3million)

If use standard/fully connected Network

~~3 billion~~ 3 billion parameters: very large



Problems: difficult to prevent overfitting

computation. Memory: Infeasible

Solution: better implement convolution operation.

更好地实现卷积运算

(CNN的基本构建块)

2. Edge Detection

Vertical Edge

Horizontal Edge

1. Vertical Edge Detector

$$\begin{bmatrix} 3 & 0 & 0 & 1 & -1 & -1 & -1 \\ 1 & 5 & 0 & 8 & -1 & -1 & -1 \\ 2 & 7 & 0 & 2 & -1 & -1 & -1 \\ 0 & 1 & 3 & 1 & 1 & -1 & -1 \\ 4 & 2 & 1 & 1 & 1 & -1 & -1 \\ 2 & 4 & 5 & 1 & 1 & -1 & -1 \end{bmatrix} * \begin{bmatrix} 1 & 0 & -1 \\ 0 & 1 & -1 \\ 1 & 0 & -1 \end{bmatrix} = 4 \times 4 \text{ Matrix}$$

6×6 (Image)

~~Another~~ Image.

Function: conv-forward

$\left\{ \begin{array}{l} \text{conv2d} \\ \text{tf.nn.conv2d} \end{array} \right.$

$10 \ 10 \ 10 \ 0 \ 0 \ 0$

$$\begin{bmatrix} \dots & \dots & \dots & 1 & 0 & -1 \\ \dots & \dots & \dots & 1 & 0 & -1 \\ \dots & \dots & \dots & 1 & 0 & -1 \\ \dots & \dots & \dots & 1 & 0 & -1 \\ \dots & \dots & \dots & 1 & 0 & -1 \end{bmatrix} = \begin{bmatrix} 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \end{bmatrix}$$

6×6

$$\begin{bmatrix} \square \end{bmatrix} * \begin{bmatrix} \square \end{bmatrix} \rightarrow \begin{bmatrix} \square \end{bmatrix}$$

(找到高暗过渡边缘)

2. Vertical and Horizontal Edge Detection

$$\begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix} \xrightarrow{\text{rotate } 90^\circ} \begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix}$$

Vertical

Horizontal

3. different Filter to detect edges

$$\begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix}$$

Sobel Filter
增加中间行
(more stable)

$$\begin{bmatrix} 3 & 0 & -3 \\ 10 & 0 & -10 \\ 3 & 0 & -3 \end{bmatrix}$$

Scharr Filter

$$\xrightarrow{\text{Bankforward 9 Parameters}} \begin{bmatrix} w_1 & w_2 & w_3 \\ w_4 & w_5 & w_6 \\ w_7 & w_8 & w_9 \end{bmatrix}$$

0°

90°

45°

70°

Any Edge

$6 \times 6 * 3 \times 3 \rightarrow 4 \times 4$

Image ^(CONV) Filter (Good) 边缘检测器

三. Padding. (对基本卷积操作的改进)

$$6 \times 6 * 3 \times 3 \rightarrow 4 \times 4$$

$$n \times n \quad f \times f \quad (n-f+1) \times (n-f+1)$$

Problems: 图像变小.

图像边缘失去大量信息.

填充像素:



$$6 \times 6 \rightarrow 8 \times 8$$

$$n \times n$$

$$* 3 \times 3 \rightarrow 6 \times 6$$

$$f \times f$$

$$(n+2p-f+1) \times (n+2p-f+1)$$

填充多少像素? Valid and Same Convolutions

$$\left\{ \begin{array}{l} \text{Valid: No padding, } n \times n * f \times f \rightarrow (n-f+1) \times (n-f+1) \\ \qquad \qquad \qquad (p=0) \quad 6 \times 6 \quad 3 \times 3 \quad 4 \times 4 \end{array} \right.$$

$$\left. \begin{array}{l} \text{Same: Output Size = Input Size } (n+2p-f+1) \times (n+2p-f+1) \\ \qquad \qquad \qquad p = \frac{f-1}{2}, (n+2p-f+1=n) \end{array} \right.$$

IV. Strided Convolutions

$$\begin{array}{l} n \times n \text{ Image} * f \times f \text{ filter} \rightarrow \lfloor \frac{n+2p-f+1}{s} \rfloor \times \lfloor \frac{n+2p-f+1}{s} \rfloor \\ \text{padding } p \\ \text{Stride } s \\ \boxed{\lfloor z \rfloor = \text{floor}(z)} \end{array}$$

分子不可整除时, 向下取整.

V. Cross-correlation 交叉相关(卷积操作).

Filter:

$$\begin{array}{|c|c|c|} \hline 3 & 4 & 5 \\ \hline 1 & 0 & 2 \\ \hline -1 & 9 & 7 \\ \hline \end{array}$$

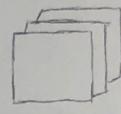
$$(A * B) * C = A * (B * C)$$

In Math: Associativity.

$$\begin{array}{|c|c|c|} \hline 7 & 2 & 5 \\ \hline 9 & 0 & 4 \\ \hline -1 & 1 & 3 \\ \hline \end{array}$$

六. Convolutions Over Volume.

1. Convolutions on RGB Images



$$6 \times 6 \times 3$$

$$* 3 \times 3 \rightarrow 4 \times 4$$

(RGB 3 通道分别有 Filter)

Height \times Width \times channel \rightarrow 27 parameters

2. Multiple Filters



(Vertical Edge)

$$= 4 \times 4$$

$3 \times 3 \times 3$

* (Horizontal Edge)



$$= 4 \times 4$$

$3 \times 3 \times 3$

$$n \times n \times n_c * f \times f \times n_c \rightarrow (n-f+1) \times (n-f+1) \times n_c$$

Filters

七. One layer of a convolutional Network

$$\begin{array}{l} 6 \times 6 \times 3 * \underbrace{3 \times 3 \times 3}_{W^{[0]}} \rightarrow 4 \times 4 + b_1 \rightarrow 4 \times 4 \\ * \underbrace{3 \times 3 \times 3}_{W^{[1]}} \rightarrow 4 \times 4 + b_2 \rightarrow 4 \times 4 \\ W^{[0]}, a^{[0]} = Z^{[1]} \end{array}$$

$$\left\{ \begin{array}{l} Z^{[1]} = W^{[1]} \cdot a^{[0]} + b^{[1]} \\ a^{[1]} = g(Z^{[1]}) \end{array} \right.$$

Number of parameters in one layer

If you have 10 filters that are $3 \times 3 \times 3$ in one layer of a neural network, how many parameters does the layer have?

$$\boxed{1}, \boxed{2}, \dots, \boxed{10} \quad \begin{matrix} \rightarrow 3 \times 3 \times 3 \\ \text{1不通过} \\ \text{拟合} \end{matrix}$$

$$27 \text{ parameters} + \text{bias} = 28 \text{ parameters}$$

$$280 \text{ parameters}$$

(无论输入图像多大, 参数数量固定是280, 而用这10个Filter来检测各种特征, 非常大的图像只是很少的参数).

If layer l is a convolution layer :

$f^{(l)}$: filter size

$p^{(l)}$: padding

$s^{(l)}$: stride.

$n_{\text{filters}}^{(l)}$: number of filters

Each Filter is : $f^{(l)} \times f^{(l)} \times n_{\text{filters}}^{(l-1)}$

Activations : $a^{(l)} \rightarrow n_H^{(l)} \times n_W^{(l)} \times n_{\text{filters}}^{(l)}$

$$A^{(l)} \rightarrow \underbrace{m \times n_H^{(l)} \times n_W^{(l)} \times n_{\text{filters}}^{(l)}}_{n_H \times n_W \times n_{\text{filters}}}$$

Weights : $f^{(l)} \times f^{(l)} \times n_{\text{filters}}^{(l-1)} \times n_{\text{filters}}^{(l)}$

Bias : $n_{\text{filters}}^{(l)} - (1, 1, 1, n_{\text{filters}}^{(l)})$

Why Convolutions ?

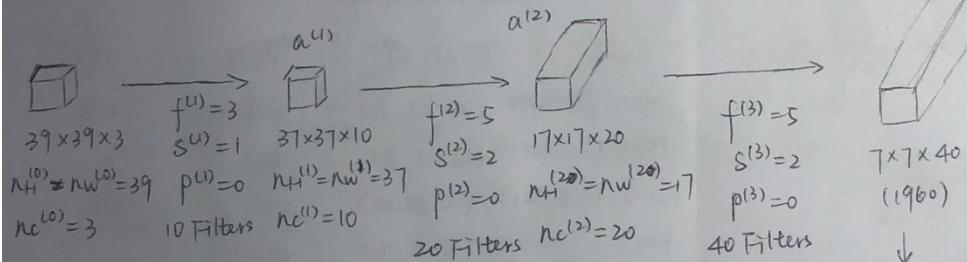
1. Parameter Sharing .

A feature detector (e.g. vertical edge detector) that's useful in one part of the image is probably useful in another part of the image .

2. Sparsity of Connections

In each layer, each output value depends only on a small number of inputs .

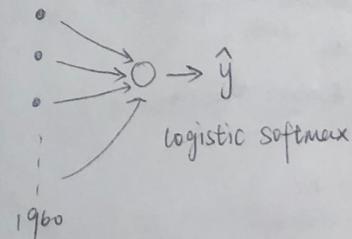
11. Simple Convolutional Network Example .



Types of layer in a convolutional Network :

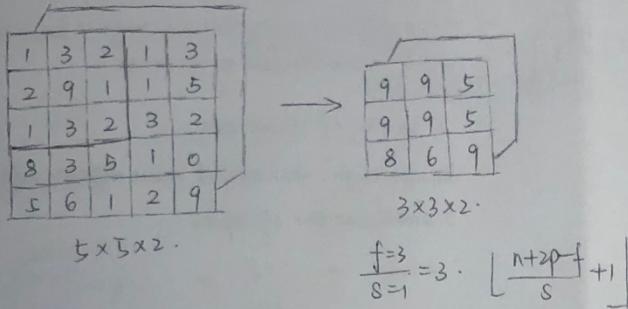
- { convolution (Conv)
- pooling (POOL)
- Fully-Connected (FC) .

(使用反向传播来训练这些网络) .

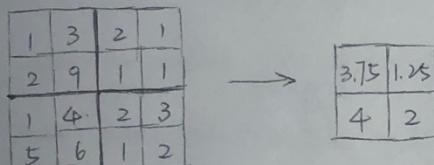


TD. Pooling Layers.

1. Max Pooling



2. Average Pooling



3. Summary of Pooling.

Hyperparameters:

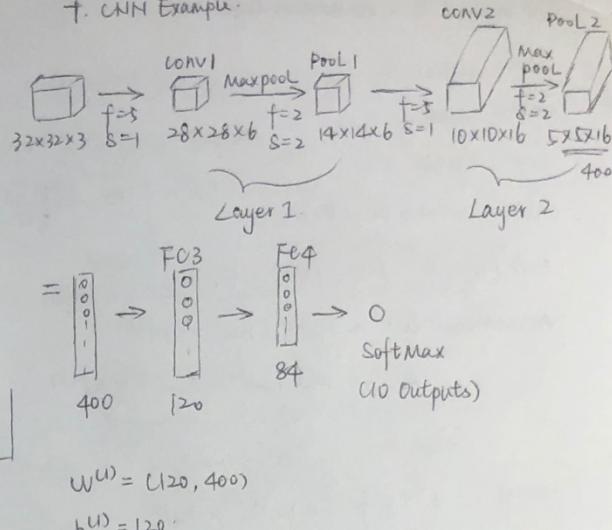
$$\begin{cases} f: \text{Filter Size.} \\ s: \text{Stride.} \end{cases} \quad \begin{cases} f=2, s=2 \\ f=3, s=2 \end{cases}$$

Max or Average Pooling.

$n_h \times n_w \times n_c$.

$$\left\lfloor \frac{n_h-f+1}{s} \right\rfloor \times \left\lfloor \frac{n_w-f+1}{s} \right\rfloor \times n_c$$

+ CNN Example.



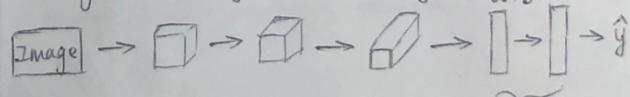
输入下去， n_h, n_w 会变小， n_c 会变大。

	Activation Shape	Activation Size	# Parameters
Input:	(32, 32, 3)	3072	0
CONV1 ($f=5, s=1$)	(28, 28, 8)	6272	208
POOL1	(14, 14, 8)	1568	0
CONV2 ($f=5, s=1$)	(10, 10, 16)	1600	416
POOL2	(5, 5, 16)	400	0
FC3	(120, 1)	120	48001
FC4	(84, 1)	84	10081
SoftMax	(10, 1)	10	841

多看些实例：如何用构件来创造高效的神经网络
如何组合。

Putting it Together :

Training Set $(x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})$ w, b.



$$\text{Cost } J = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

use gradient descent to optimize parameters to reduce J

Ch 2 Deep Convolutional Models (Case studies)

Outline:

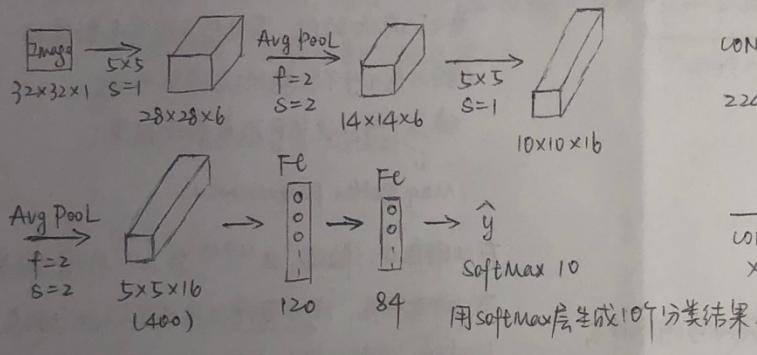
- Classic Networks:
 - LeNet-5
 - AlexNet
 - VGG

• ResNet

• Inception

1. Classic Networks

1. LeNet-5. (小型神经网络)

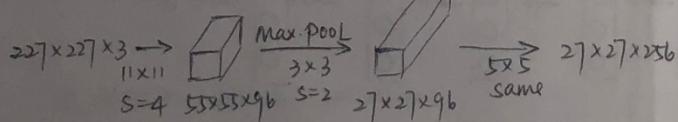


NH, NW ↓, NC ↑

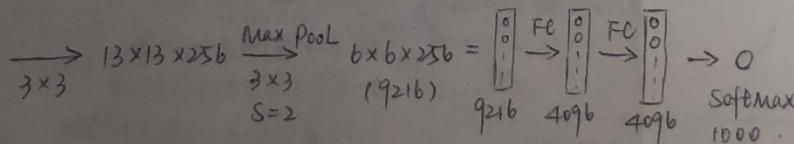
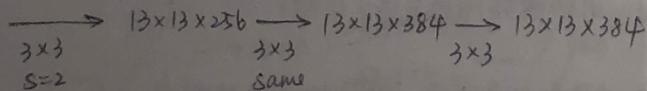
CONV, POOL, CONV, POOL, FC, FC, output.

为了节省计算量和参数数量：不同的Filter处理不同channel。
这种方法很复杂，现在不再使用。

2. AlexNet (有点像 LeNet-5, 但更大~)



MAX POOL



① AlexNet: Similar to LeNet-5, but much Larger.

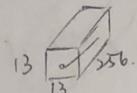
有更多隐藏神经元，在更多数据上训练。
⇒ better performance.

② ReLU 激活函数

不 \rightarrow ③ Multiple GPU.

太 \rightarrow ④ Local Response Normalization (LRN).

图 局部响应归一化：一个位置的高和宽，然后看跨
越通道的这个位置，然后归一化。
(针对处理每个位置.)



3. VGG-16. (非常大的神经网络，很简洁)

CONV = 3×3 Filter, $s=1$, same.

$224 \times 224 \times 3 \rightarrow 224 \times 224 \times 64 \rightarrow 112 \times 112 \times 64$
CONV 64 POOL $\times 2$

$\rightarrow 112 \times 112 \times 128 \rightarrow 56 \times 56 \times 128$
CONV 128 POOL $\times 2$

$\rightarrow 56 \times 56 \times 256 \rightarrow 28 \times 28 \times 256$
CONV 256 POOL $\times 3$

$\rightarrow 28 \times 28 \times 512 \rightarrow 14 \times 14 \times 512$
CONV 512 POOL $\times 3$

$\rightarrow 14 \times 14 \times 512 \rightarrow 7 \times 7 \times 512$
CONV 512 POOL $\times 3$

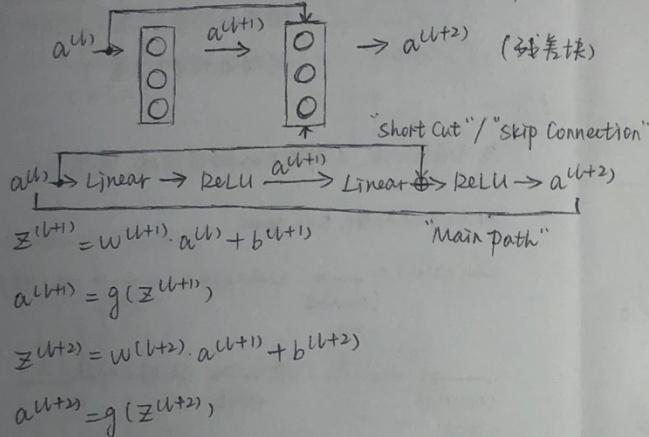
\rightarrow FC \rightarrow FC \rightarrow SoftMax
4096 4096 1000

NH, NW ↓, NC ↑. 这种操作减少或增加
很有道理

二. Residual Networks (ResNets)

太深的NN训练起来很难，因为“梯度消失”“爆炸”
 ↓
 Skip connection: 从一层中得到激活 → 突进，送给下一层 → deeper layer
 ↓
 训练很深的残差网络 (ResNets)

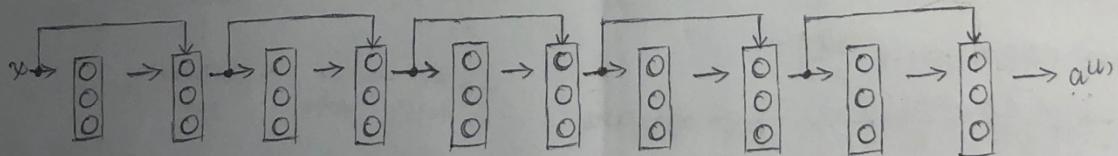
1. Residual Block.



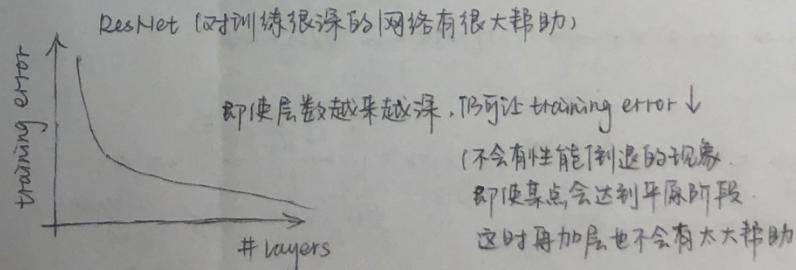
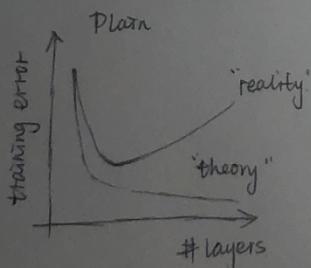
Skip connection: 几乎把两层信息传递到更深的神经网络中去。
 使用残差块可以训练更深的神经网络。

建立一个ResNet的方法：通过大量的残差块，将它们堆叠起来。
 形成一个深层网络 (Deep Network)

2. Residual Network.



5个残差块堆积在一起。⇒ 形成一个残差网络。



三. Why ResNets Work

$$\begin{aligned} x &\rightarrow \text{Big NN} \rightarrow a^{l+1} \\ x &\rightarrow \text{Big NN} \xrightarrow{a^{l+1}} \boxed{\text{ReLU } (a>0)} \rightarrow a^{l+2} \\ &\quad \text{ReLU } (a>0) \quad \uparrow \\ a^{l+2} &= g(z^{l+2}) = g(w^{l+2} \cdot a^{l+1} + b^{l+2} + a^{l+1}) \\ &= g(w^{l+2} \cdot a^{l+1} + b^{l+2}) \quad \uparrow \\ &= g(a^{l+1}) = a^{l+1} \end{aligned}$$

中间额外添加两层，所有项都学到有意义的东西的话，要比学习恒等函数要好。
 (额外外层几乎不会影响具体的表达)
 梯度下降从这里可改善最终结果。
 ↓
 May Better Performance.

在训练深度中，假设 z^{l+2} 与 a^{l+1} 在同一维度。
 有一堆卷积层，偶尔有池化层或类似池化的层。
 可用 WS 调整维度。一个充分连接的层来最终用 softmax 做预测。

$z^{l+2} + a^{l+1}$. 3x3 same. pool L. WS
 ↓
 FC. Softmax 1000.

① Network In Network and 1×1 convolutions.

1. Simple:

$$6 \times b \times 1 \cdot * [2] = 6 \times b \times 1 \\ (\text{just Multiply by } 2)$$

2.

$$6 \times b \times 32 \cdot * 1 \times 1 \times 32 = b \times b \times \# \text{Filters.} \\ (\text{ReLU})$$

将 Left $32 \times$ Filter 32 (内核).

再将一个非线性映射 ReLU 作用于它.

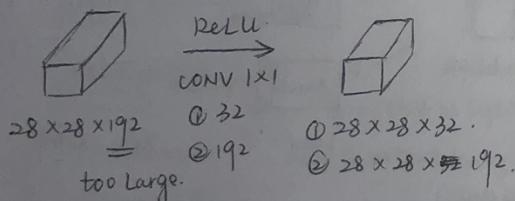
类似于：已有一个神经元，接收一个 32 个数的输入向量.

Output: 过滤器数目输出值.

↓

1×1 conv (Network In Network)

3. Ne↓.



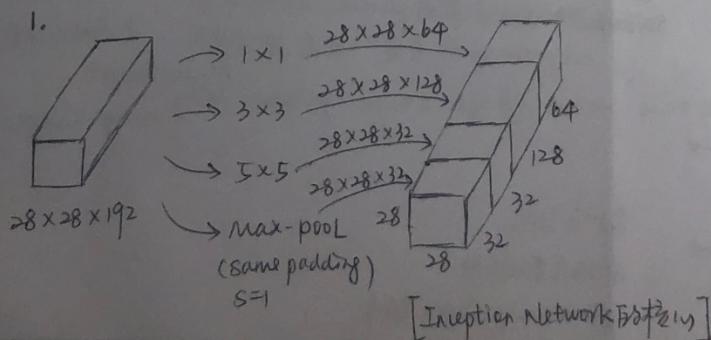
① CONV 1×1 : decrease # nc \Rightarrow 减小计算量.

pool layer: decrease nh, nw.

② CONV 1×1 : 增加非线性.

二. Inception Network Motivation.

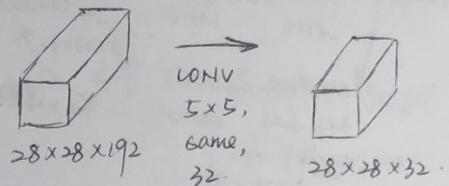
1.



Inception Network 基础特点: ① 不用只选一个卷积核的大小或 pooling.

② 所有可能都做, 将所有输出结果都连起来.

2. Problem of Computational Cost

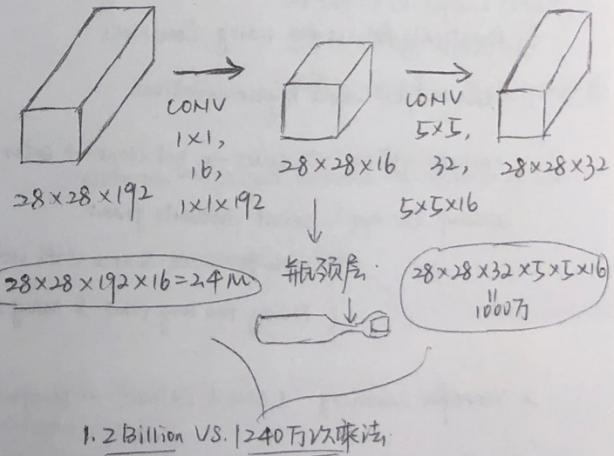


$$28 \times 28 \times 32 \times 5 \times 5 \times 192 = 120 \text{ million} \\ (1.2 \text{ billion}).$$

Too Large Comp~Cost.

↓

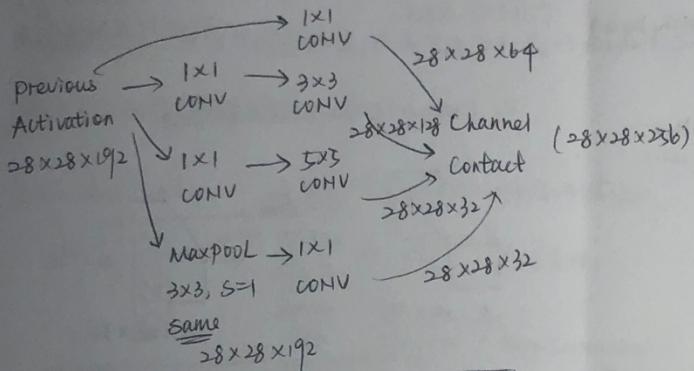
1×1 conv: decrease the Computational Cost $\frac{1}{10}$ ~



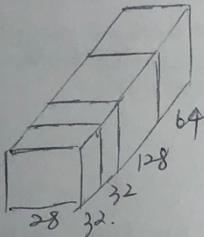
(合理实现瓶颈层, 又不影响整体性能).

六. Inception Network. (重复多次使用 Inception Module).

Inception Module :



旁枝：FC. 将隐藏层作为输入
来预测 (Normalization)
 \downarrow
prevent Overfitting.



f. Practical Advice for Using ConvNets

1. Using Open-Source Implementations

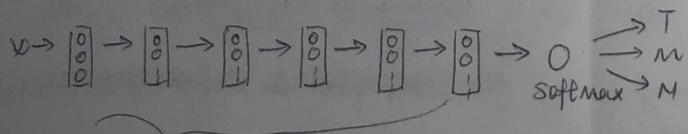
resnets github → codes → git clone + enter.

develop CV App: select favorite frame

look for open-source with codes

(Training too long time. & Many GPU. Dataset. → Transfer Learning Directly).

2. Transfer Learning (Small Dataset. → Useful -)



freeze.

freeze = 1

training parameter = 0.

freeze.

train.

{ Gradient.
Delete.

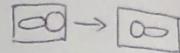
参数层越多

冻结层参数越少

自己训练 Layer 越多

3. Data Augmentation. (Adjust hyperpara.)

(1) Common Augmentation Method : ① Mirroring (垂直镜)

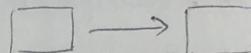


② random cropping. (随机剪裁)

③ 1-rotation Shearing (旋转扭曲)

④ Local warping (局部弯曲)

(2) Color Shifting



R.G.B.

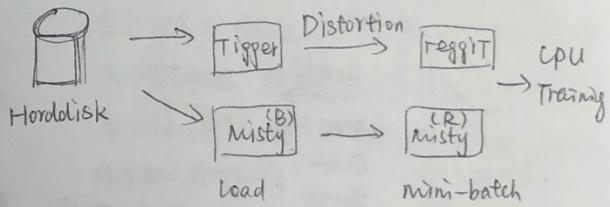
+20, -20, +20. (不同的扰动)

-20, +20, +20

+5, 0, +50

(3) Implement Distortions during Training.

(CPU Thread.)



load

mini-batch

Using Dataset. = Initialize Network.

(Weights)

↓ Gradient.

Renew weights.

& Network Layers.

Ch 3. Detection Algorithms.

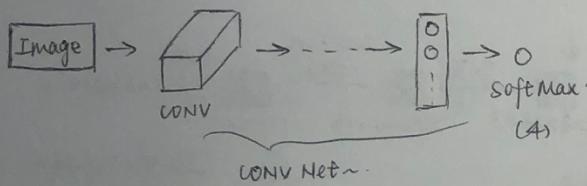
1. Object Localization.

1. Localization and Detection.

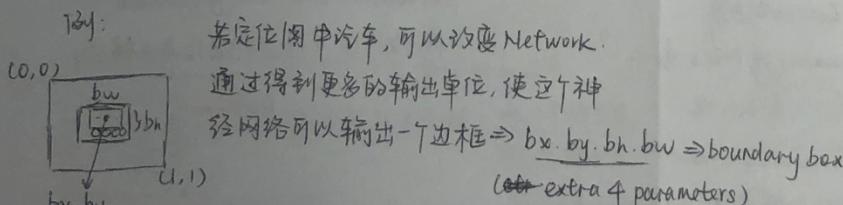
- Image Classification
- Classification with Localization } [Commonly] Object. (识别正中央的~~一个明显对象~~)
- Detection \Rightarrow Multiple Objects (识别来自不同类别的多个对象)

2. Classification with Localization

在多层卷积神经网络中输入一张图片，这会教导我们的卷积网络中有一个特征向量，将特征向量送给一个 softmax，输出预测结果。



- By:
 分类识别
 ① pedestrian
 ② car \Rightarrow 有 4 个结果的 SoftMax.
 ③ motorcycle
 ④ background. Need to output b_x, b_y, b_h, b_w
 Class Label (4)



用 supervised Learning 让算法输出的不仅仅是
是坐标，还有 4 个变量 \Rightarrow boundary box.

(Estimation: $b_x=0.5, b_y=0.7, b_h=0.3, b_w=0.4$)

Defining the Target Label y :

$$y = \begin{bmatrix} p_C \\ b_x \\ b_y \\ b_h \\ b_w \\ C_1 \\ C_2 \\ C_3 \end{bmatrix}$$

损失值是元素平方差之和

$$y = \boxed{\text{0000}}$$

$$\mathcal{L}(y, \hat{y}) = \begin{cases} (\hat{y}_1 - y_1)^2 + \dots + (\hat{y}_8 - y_8)^2 & y_1 = 1 \\ (\hat{y}_1 - y_1)^2 & y_1 = 0 \end{cases}$$

其它元素都不重叠。
只关注 Network.

$y_1=1$ if. pay attention to Network \Rightarrow 可用均方差 penalize
 $y_1=0$ if. don't care other elements. 预测平方偏差，并从 8 个组件的数据输出

& 也可用 Log likelihood 特征损失
 对于 C_1, C_2, C_3 . softmax output,
 其中之一的元素，可用均方差
 如果 p_C 可用类似逻辑回归损失.

Network \rightarrow Output numbers \rightarrow Location of sth

powerful Idea

二. Landmark Detection.

有标签的训练集：

$$x, y \left\{ \begin{array}{l} \text{f}_1, x, f_2, y \\ f_2, x, f_3, y \\ f_3, x, f_4, y \\ f_4, x, f_5, y \\ \vdots \\ f_{64}, x, f_{64}, y \end{array} \right.$$

想要识别的
特征的 (x, y)
标签必须在不
同图片中保持
一致性。

三. Object Detection. (ConvNet. Sliding-Wins)

1. car detection example.

Training Set.:

x	y
car	1
car	1
car	1
Tree	0
Mountain	0

2. Sliding Windows detection.

滑动小窗口遍历Image. \Rightarrow ConvNet.
(stride).

将小窗口Area \uparrow . 重叠遍历Image.

Algorithm: Input: Small Frame. \rightarrow ConvNet.
Output: $y=1 \rightarrow$ Detect a ~~car~~ car

Disadvantage: Large Computation Cost.

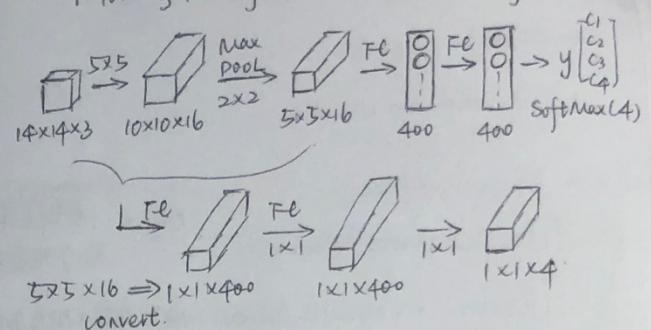
每个图像都单独通过ConvNet运算
① 如果 stride↑ 不密集. \rightarrow Number of Frame ↓

\rightarrow 粗糙的颗粒度会影响算法表现.

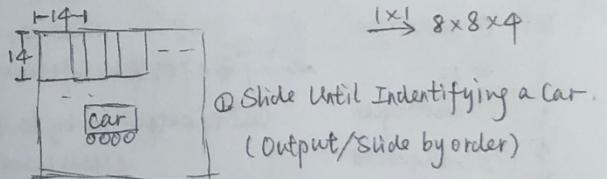
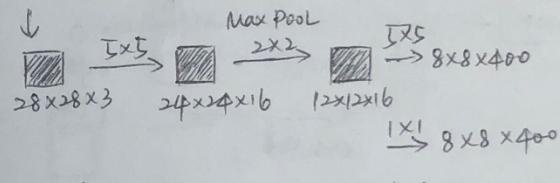
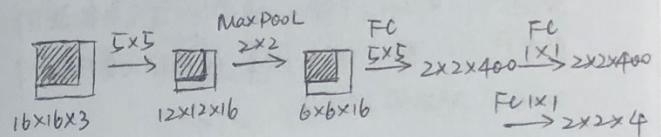
② too small stride. \rightarrow Number of Frame ↑
 \rightarrow 精细的颗粒度 \rightarrow 小区域密集 \rightarrow Large Loss

IV. Convolutional Implementation of Sliding Windows

1. Turning FC Layer into Convolutional Layers.



2. Convolution Implementation of Sliding Windows.



① Slide Until Identifying a Car.
(Output/Slide by order)

✓ ② ConvNet \rightarrow Forward propagation.
(Better) \rightarrow y predict.

高效许多. (但边框定位并不精确)

Problem.

3. Bounding Box Predictions.

(Solution for Sliding-Wins Algorithms.)

More accurate to predict Boundary Frame)

Car 不一定是“方框”。

{ 方框 不一定完全准确框住车 }

1. YOLO Algorithm (检测较困难)

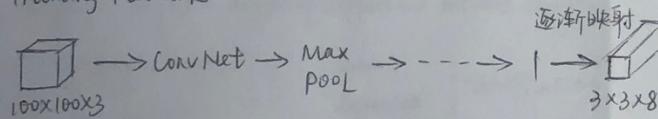
Labels for Training:

For each grid cell:

$$y: \begin{bmatrix} PC \\ bx \\ by \\ bw \\ bh \\ C_1 \\ C_2 \\ C_3 \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ bx \\ by \\ bw \\ bh \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

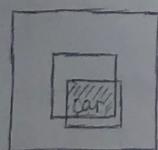
$$9 = 3 \times 3 \quad (\text{8 Dimension } y) \quad (3 \times 3 \times 8)$$

Training Network:



7. Intersection Over Union.

1. Evaluating Object Localization.



The result is good or bad?
交并比或 IoU 函数: $\frac{\text{两边界框交集}}{\text{两边界框并集}}$

$$= \frac{\text{Size of } \square}{\text{Size of } \square}$$

"Correct" if IoU ≥ 0.5

More Generally, IoU is a measure of the overlap between two bounding boxes. IoU 是两个边界框重叠程度的一个度量值。

4. Non-Max Suppression. 非最大值抑制.

Algorithm may have many testings for the same target.

Non-Max Suppression makes sure the algorithm has only one testing for each target (Object).

运行算法时，对每个目标，也许会得到多个检测结果

清理检测，每辆车只会得到一个检测结果。

④ The probability of each testing result (P_c, C_i)

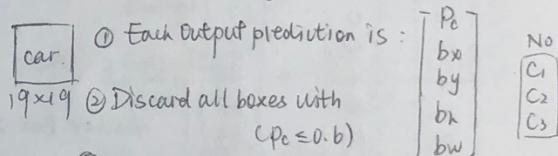
Select the largest

& Suppress other results.

1. Non-Max: 输出有着最大可能性的分类判断

抑制那些 Non-Max 的邻近方程。

2. Non-Max Suppression Algorithm

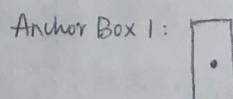


③ Remaining Boxes:

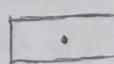
- Pick the box with the largest P_c output that as a prediction.
- Discard any remaining box with $IoU \geq 0.5$ with the box output is the previous step.

11. Anchor Box.

1. Overlapping Objects



Anchor Box 1:

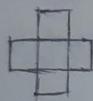


$$y = \begin{bmatrix} P_0 \\ b_x \\ b_y \\ b_w \\ b_h \\ C_1 \\ C_2 \\ C_3 \\ P_c \\ C_3 \end{bmatrix} \left\{ \begin{array}{l} \text{Anchor Box 1 (8 parameters)} \\ \text{Anchor Box 2 (8 parameters)} \end{array} \right.$$

2. Anchor Box Algorithm (处理两个对象出现在同一格子的情况)

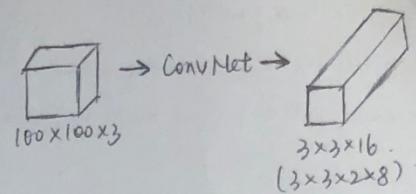
Previously: Each object in training image is assigned to grid cell that contains that object's midpoint.
将训练图像上的对象指定给对应目标中心点的网格单元。 (Output y: $3 \times 3 \times 8$)

With Two Anchor Boxes: Each object in training image is assigned to grid cell that contains object's midpoint and anchor box for the grid cell with highest IoU.
 $y: 3 \times 3 \times 16$



$$y = \left[\begin{array}{c|c|c} p_c & 1 & 0 \\ \hline bx & bx & ? \\ by & by & ? \\ bh & bh & ? \\ bw & bw & ? \\ \hline C_1 & 1 & ? \\ C_2 & 0 & ? \\ C_3 & 0 & ? \\ \hline p_c & 1 & 1 \\ bx & bx & bx \\ by & by & by \\ bh & bh & bh \\ bw & bw & bw \\ \hline C_1 & 0 & 0 \\ C_2 & 1 & 1 \\ C_3 & 0 & 0 \end{array} \right] \quad \text{Anchor Box 1}$$

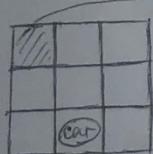
$$\left. \begin{array}{c|c|c} p_c & 1 & 0 \\ \hline bx & bx & ? \\ by & by & ? \\ bh & bh & ? \\ bw & bw & ? \\ \hline C_1 & 0 & ? \\ C_2 & 1 & ? \\ C_3 & 0 & ? \\ \hline p_c & 1 & 1 \\ bx & bx & bx \\ by & by & by \\ bh & bh & bh \\ bw & bw & bw \\ \hline C_1 & 0 & 0 \\ C_2 & 1 & 1 \\ C_3 & 0 & 0 \end{array} \right] \quad \text{Anchor Box 2.}$$



Outputting the non-max suppressed outputs:

- ① For each grid cell, get 2 predicted bounding boxes
- ② Get rid of low probability predictions.
- ③ For each class (pedestrian, car, motorcycle)
Use non-max suppression to generate final predictions.

3. Putting it together: YOLO Algorithm



{ Pedestrian
car
motorcycle. }

① y is $3 \times 3 \times 2 \times 8$
 \downarrow
 $(\text{维度}) \cdot 5 + 3C$.
 Anchor Box.

$3 \times 3 \times 16$.

② $19 \times 19 \times 16$
 $19 \times 19 \times 5 \times 8 \rightarrow 19 \times 19 \times 40$
 \downarrow
 5 Anchor Boxes.

$$y = \left[\begin{array}{c|c|c} p_c & 1 & 0 \\ \hline bx & ? & ? \\ by & ? & ? \\ bh & ? & ? \\ bw & ? & ? \\ \hline C_1 & ? & ? \\ C_2 & ? & ? \\ C_3 & ? & ? \\ \hline p_c & 0 & 1 \\ bx & ? & ? \\ by & ? & ? \\ bh & ? & ? \\ bw & ? & ? \\ \hline C_1 & ? & ? \\ C_2 & ? & ? \\ C_3 & ? & ? \\ \hline p_c & 1 & 0 \\ bx & bx & by \\ by & by & bh \\ bh & bh & bw \\ bw & bw & 0 \\ \hline C_1 & 0 & 1 \\ C_2 & 1 & 0 \\ C_3 & 0 & 0 \end{array} \right]$$

Ch4 Face Recognition.

- Face verification VS Face Recognition.

1. Verification.

- Input Image, name / ID

- Output whether the input image is that of the claimed person.

2. Recognition.

- Has a database of K persons.
- Get an input image.
- Output ID if the image is any of the K persons
(for "not recognised")

二. One-Shot Learning. 单样本学习 (输入2张图, 看“相似”or“不同”)

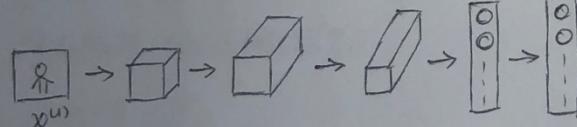
1. Learning from one example to recognize the person again

2. Learning a "similarity" Function.

$d(\text{img}_1, \text{img}_2) = \text{degree of difference between images.}$

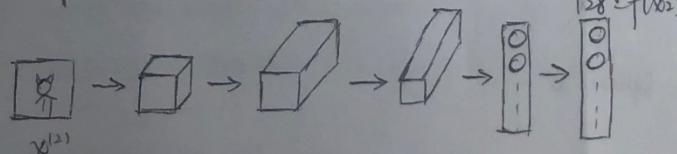
If $d(\text{img}_1, \text{img}_2) \leq T$ "same"
 $> T$ "different" } Verification

三. Siamese Network.



专注向量本身, "128". 从 Network 深处的某个全连接层计算而来。

$128: f(x^{(1)}) = \text{encoding of } x^{(1)}$



使用相同的参数, 得到一个不同的, 128个数字组成的向量。

$$d(x^{(1)}, x^{(2)}) = \|f(x^{(1)}) - f(x^{(2)})\|_2^2$$

Parameters of NN define an encoding $f(x^{(i)})$

Learn parameters so that:

- If $x^{(i)}, x^{(j)}$ are the same person, $\|f(x^{(i)}) - f(x^{(j)})\|^2$ is small

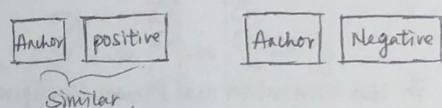
- If $x^{(i)}, x^{(j)}$ are different, $\|f(x^{(i)}) - f(x^{(j)})\|^2$ is large.

当设置所有层的神经网络的参数, 可以使用
反向传播来使参数以满足条件。

IV. Triplet Loss Function (梯度下降的三元组损失函数)

1. Aim: 为了学习神经网络的参数, 并以此获得
一个优秀的人脸图片的编码。

2. Learning Objective:



$$\text{Want: } \frac{\|f(A) - f(P)\|^2}{d(A, P)} \leq \frac{\|f(A) - f(N)\|^2}{d(A, N)}$$

$$\frac{\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2}{d(A, P) - d(A, N)} \leq 0$$

+ α
(margin para -)

Given 3 Images
A, P, N:

$$L(A, P, N) = \max(\underbrace{\|f(A) - f(P)\|^2}_{>0} - \underbrace{\|f(A) - f(N)\|^2 + \alpha}_{>0}, 0)$$

$$J = \sum_{i=1}^m L(A^{(i)}, P^{(i)}, N^{(i)})$$

Training set: DK pairs of $1K$ persons

如果每个人只有一张 pic, 无法训练 (许多对)

A, P, N
同一人的

3. Choosing the triplets A, P, N

During training if A, P, N are chosen randomly,

$d(A, P) + \alpha \leq d(A, N)$ is easily satisfied.

$$\|f(A) - f(P)\|^2 + \alpha \leq \|f(A) - f(N)\|^2$$

Choose triplets that're "hard" to train on:

$$d(A, P) + \alpha \leq d(A, N)$$

$$d(A, P) \approx d(A, N)$$

选择三元组，增强学习算法的计算效率。

4. Training Set Using Triplet Loss:

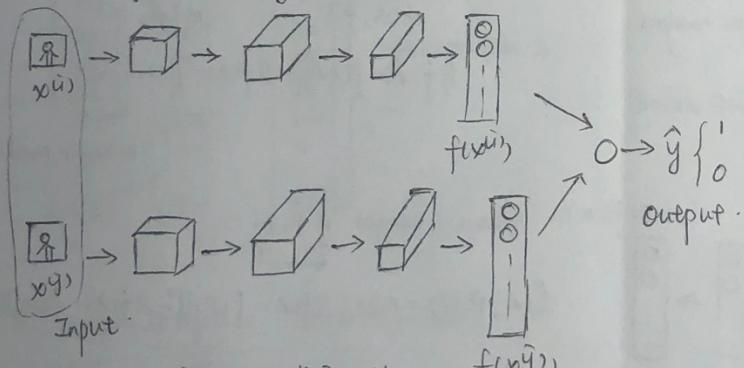
$$J_{\text{min.}}(d(x^i), d(y))$$

2. Face Verification Supervised Learning

x	y	
□	1	"same"
□	0	"different"
:	0	
Image	1	

五. Face Verification and Binary Classification.

1. Learning the Similarity Function



不是将 Data 直接输入，而是要计算
编码之间的不同。

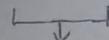
$$\begin{aligned} \hat{y} &= \sigma \left(\sum_{k=1}^{28} \|f(x^i)_k - f(x^j)_k\|^2 \times w_k + b \right) \\ &\approx \frac{(f(x^i)_k - f(x^j)_k)^2}{f(x^i)_k + f(x^j)_k} \quad (\text{Chi-Square } \chi^2) \end{aligned}$$

Precompute: 无需每次都为每个数据库中的员工计算编码。

Ch.5. Neural Style Transfer.

1. Neural Transfer.

Content. Style.

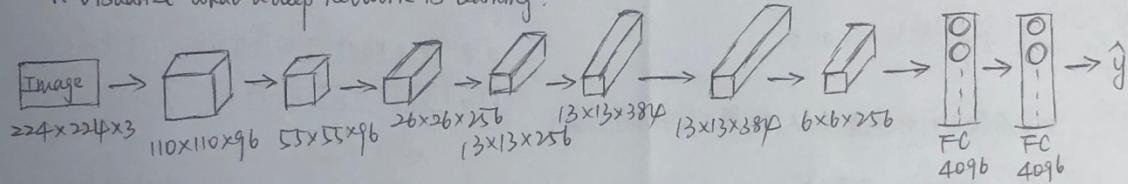


Generated Image.

In order to achieve neural transfer, look at features extracted by ConvNet at various layers. the shallow and the deeper layers of a ConvNet.

2. Deep ConvNets Learning.

1. Visualize what a deep network is learning.



⇒ Pick a unit in layer 1. Find the nine image patches that maximize the unit's activation.

⇒ Repeat for other units.

2. Visualize deep layers: Layer 2, 3, ...

层数越深，模式或特征越复杂

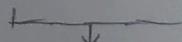
Layer 1: 边.

Layer 2: 文字.

Deeper Layer: very complex objects

3. Neural Style Transfer Cost Function.

Content (C) Style (S)



Generated Image (G).

$$J(G) = \alpha \cdot J_{content}(C, G) + \beta \cdot J_{style}(S, G)$$

Find the generated image G:

① Initiate G randomly: $G \text{ (} 100 \times 100 \times 3 \text{ - RGB)}$.

② Use gradient descent to minimize $J(G)$.

$$G = G - \frac{\partial J(G)}{\partial G} \quad (\text{相当于有更新图像 } G \text{ 的像素值})$$

IV. Content Cost Function.

$$J(G) = \alpha \cdot J_{content}(C, G) + \beta \cdot J_{style}(S, G)$$

① Use hidden layer l to compute $J_{content}$.

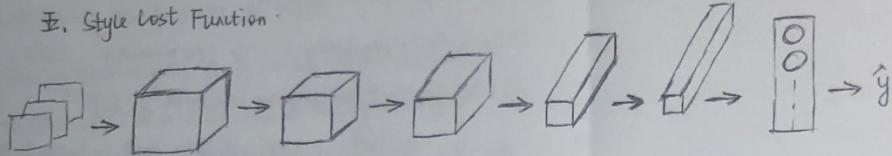
② Use pre-trained ConvNet (E.g. VGG16 network ...)

③ Let $a^{l+1}(C)$ and $a^{l+1}(G)$ be the activation of layer l on the images.

④ If $a^{l+1}(C)$ and $a^{l+1}(G)$ are similar, both images have similar content.

$$J_{content}(C, G) = \frac{1}{2} \| a^{l+1}(C) - a^{l+1}(G) \|^2$$

五. Style Loss Function.



Use Layer l 's activation to measure "style".

Define style as correlation between activations across channels.

Style Matrix.

$$a_{i,j,k}^{[l]} = \text{activation at } (i, j, k). \quad G^{[l]} \text{ is } n_h^{[l]} \times n_w^{[l]} \times n_c^{[l]}$$

$\downarrow \quad \downarrow \quad \downarrow$
 $H \quad W \quad C$

讲述每对通道的相关性.

$$\Rightarrow G_{kk'}^{[l](s)} = \sum_{i=1}^{n_h^{[l]}} \sum_{j=1}^{n_w^{[l]}} a_{ijk}^{[l]} \cdot a_{ijk'}^{[l]}$$

$k=1 \dots n_c^{[l]}$, k 要遍历不同的通道。(从1到第 l 层的总通道数)

$$\Rightarrow G_{kk'}^{[l](G)} = \sum_{i=1}^{n_h^{[l]}} \sum_{j=1}^{n_w^{[l]}} a_{ijk}^{[l]} \cdot a_{ijk'}^{[l]} \quad \text{"Gram matrix"}$$

$$\begin{aligned} J_{\text{style}}(S, G) &= \frac{1}{(\dots)} \| G^{[l](s)} - G^{[l](G)} \|_F^2 \\ &= \frac{1}{(2n_w n_h n_c)^2} \cdot \sum_{k,k'} (G_{kk'}^{[l](s)} - G_{kk'}^{[l](G)})^2 \end{aligned}$$

$$J_{\text{style}}(S, G) = \sum_l J_{\text{style}}^{[l]}(S, G)$$

↓

$$J(G) = \alpha J_{\text{content}}(C, G) + \beta J_{\text{style}}(S, G)$$

总代价函数 = 内容代价函数 + 风格代价函数,

然后使用梯度下降法, 或采用更复杂的优化算法

来试图寻找 Image $G \Rightarrow J(G) \min. \Rightarrow$ 神经网络艺术创作.

①. 1D, 3D generalizations
of models.

②. Convolutions in 2D and 1D. (2D \rightarrow 1D).

$\boxed{\text{1111}}$	$*$	$\boxed{\text{111}}$	$14 \times 14 \times 3 * 5 \times 5 \times 3$
		\downarrow	
		$2D \text{ Input Image}$	$2D \text{ Filter}$
		14×14	5×5
		$10 \times 10 \times 16$	
		$5 \times 5 \times 16$	
		$6 \times 6 \times 32$	

By:

11张图。(按时间顺序显示每一帧的图像)

$\boxed{1 \ 20 \ 15 \ 3 \ 18 \ 12 \ 4 \ 17} \quad (14)$

$*$	$\boxed{1 \ 3 \ 10 \ 3 \ 1}$	$14 \times 1 * 5 \times 1$
		\downarrow
		$(10 \times 1) * 5 \times 16$
		6×32

用同样的一组特征去检测在这些时间序列中不同位置的11帧跳。⇒ ConvNet 可用于 1D data. (递归神经网络).

二. 3D convolution.

$\boxed{\text{1111}}$	$*$	$\boxed{\text{111}}$
		$3D \text{ volume}$
		$3D \text{ filter}$

$$14 \times 14 \times 14 \times 1 * 5 \times 5 \times 5 \times 1$$

$$\Rightarrow 10 \times 10 \times 10 \times 16 * 5 \times 5 \times 5 \times 16 \quad (16 \text{ filters}).$$

$$\Rightarrow 6 \times 6 \times 6 \times 32. \quad (32 \text{ filters})$$

Name: Xiaofan Wang

Sequence Models.

Chapter 1. Recurrent Neural Network.

1. Sequence Models.

1. Example of Sequence Data.

Speech Recognition.

Music Generation

Sentiment Classification.

DNA Sequence Analysis

Machine Translation

Video Activity Recognition.

Name Entity Recognition.

2. Notation.

x : Harry Potter and Hermione Granger invented a new spell.

$x^{(1)}, x^{(2)}, x^{(3)} \dots x^{(T)}$ - 上标代表单词位置

$$\Rightarrow y: 1 \ 1 \ 0 \ 1 \ 1 \ 0 \ 0 \ 0 \ 0 \quad \left\{ \begin{array}{l} Tx = 9 \text{ (句子长度)} \\ Ty = 9 \end{array} \right.$$

$y^{(1)}, y^{(2)}, \dots, y^{(T)}, \dots, y^{(9)}$

$x^{(i)(t)}$: T 的序列中的第 i 个元素.

$$\left\{ \begin{array}{l} Tx^{(i)} = 9 \text{ (句子长度)} \\ Ty^{(i)} = 9 \end{array} \right.$$

3. Representing Words

x : Harry Potter and Hermione Granger invented a new spell.

$x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(T)}, \dots, x^{(9)}$

Vocabulary.

将每个词作为一个“one-hot”的表示

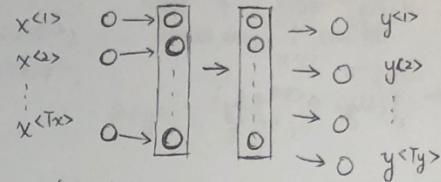
a	1
Aaron	2
;	3
And	367
;	4
Harry	4075
;	5
Potter	6830
;	6
Zulu	10000 (将1万词作为拟构建字典)

$$x^{(1)}: \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow 4075$$

$$x^{(2)}: \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow 6830$$

2. Recurrent Neural Network Model.

1. Why not a standard network.



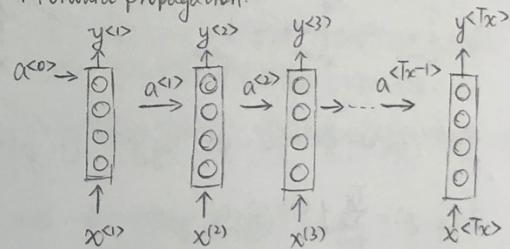
将 9 个 Input Words 放入标准的神经网络中。
然后输出 9 个 Output Words. $\left| \begin{array}{l} 0 \Rightarrow \text{每个单词是否是人名} \\ 1 \Rightarrow \text{的一部分} \end{array} \right.$

Problems:

① Inputs, Outputs can be different lengths in different examples.

② Don't share features learned across different positions of text. 像这种朴素的 NN，并不会共享那些从不同文本位置学到的特征。

2. Forward Propagation.



$$\begin{aligned} a^{(0)} &= \vec{0} & a^{(1)} &= g_1(W_a a^{(0)} + W_x \cdot x^{(1)} + b_a) \leftarrow \text{tanh/ReLU} \\ y^{(1)} &= g_2(W_y a^{(1)} + b_y) \leftarrow \text{sigmoid} \\ a^{(2)} &= g(W_a a^{(1)} + W_x \cdot x^{(2)} + b_a) \\ y^{(2)} &= g(W_y a^{(2)} + b_y) \end{aligned}$$

3. Simplified RNN Notation.

$$a^{<t>} = g(W_{aa} \cdot a^{<t-1>} + W_{ax} \cdot x^{<t>} + b_a)$$

↑ 100 ↑ 10000
(100, 100) (100, 10000)

$$\hat{y}^{<t>} = g(W_y \cdot a^{<t>} + b_y)$$

$$\Downarrow$$

$$\hat{y}^{<t>} = g(W_y \cdot a^{<t>} + b_y)$$

这里的 W_a, b_a 用于计算的参数如，激活输出的值

$$a^{<t>} = g(W_a [a^{<t-1>}, x^{<t>}] + b_a)$$

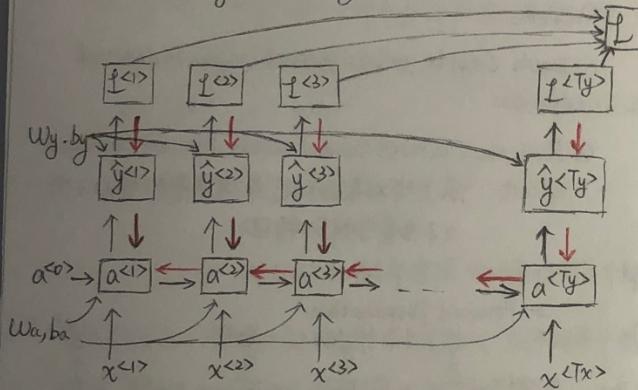
$$W_a = [W_{aa} | W_{ax}] \uparrow 100$$

(100, 10000) 100 10000

$$[a^{<t-1>}, x^{<t>}] = \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} \uparrow 100 \uparrow 10000 \uparrow 10100$$

$$[W_{aa} | W_{ax}] \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} = W_{aa} \cdot a^{<t-1>} + W_{ax} \cdot x^{<t>}$$

4. Backpropagation through time.



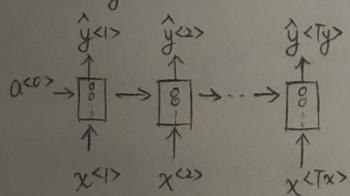
$$l^{<t>}(\hat{y}^{<t>}, y^{<t>}) = -y^{<t>} \log \hat{y}^{<t>} - (1-y^{<t>}) \log (1-\hat{y}^{<t>})$$

$$L(\hat{y}, y) = \sum_{t=1}^{Ty} l^{<t>}(\hat{y}^{<t>}, y^{<t>})$$

5. Different types of RNN:

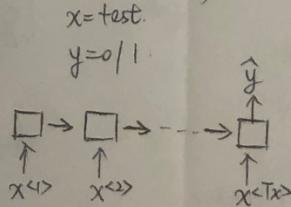
Examples of RNN Architectures:

① $Tx = Ty$.



Many-to-Many.

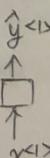
② Sentiment Classification



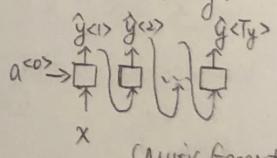
Many-To-One.

(输入: 许多单词)
(输出: 一个数字)

③ One-to-One.

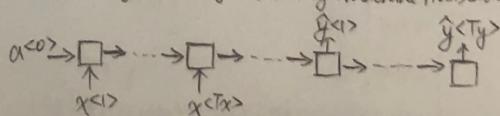


④ One To Many.



(Music Generation)

⑤ Many To Many. ($Tx \neq Ty$. Machine Translation)



三. Language Model and Sequence Generation.

1. Language Modeling.

Speech Recognition: $\begin{cases} \text{The apple and } \underline{\text{pear}} \text{ salad.} & (P_1 = 3.2 \times 10^{-13}) \\ \text{The apple and } \underline{\text{pear}} \text{ salad. } \checkmark & (P_2 = 5.7 \times 10^{-10}) \end{cases}$

语音识别系统选出第=个句子。

由于语言模型 (Language Model).

给出的两个句子的概率.

2. Language Modelling with an RNN.

Training Set: Large corpus of English text. 单词大文集.

Bj: cats average 15 hours of sleep a day. <EOS>

$y^{<1>} y^{<2>} \dots y^{<9>} \quad$ 句子结束标记.

The Egyptian Man is a bread of cat. <EOS>

<UNK> 未知单词

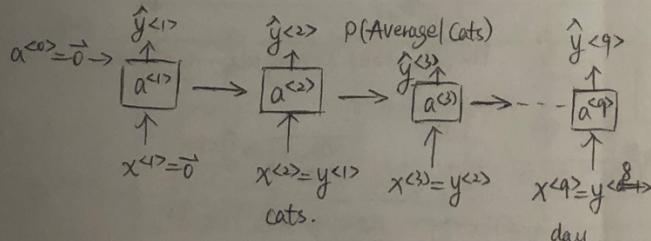
对UNK的根序建模. 而不是对特定未知词建模.

完成标记化: 将输入句子映射成各个标记. 或单词表中单词集合.

然后. 构建RNN来为这些不同序列的根序建模.

3. RNN Model (从左至右每次预测下一个单词的概率).

$p(a), p(Aaron) \dots p(zulu) p(<\text{UNK}>) p(<\text{EOS}>)$.



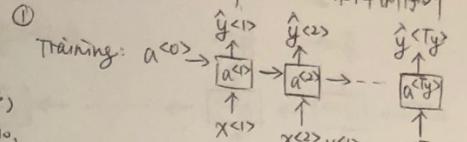
$$L(\hat{y}^{<t>}, y^{<t>}) = -\sum_i y_i^{<t>} \log \hat{y}_i^{<t>}$$

为了训练这个神经网络. 要定义一个代价函数.

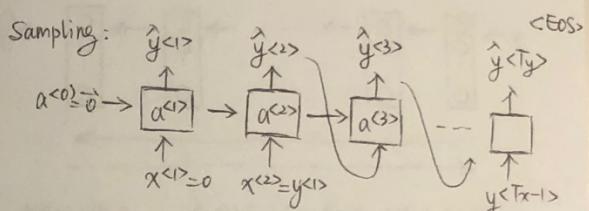
$$L = \sum_t L^{(t)}(\hat{y}^{<t>}, y^{<t>})$$

整个句子的概率分布: $p(y^{<1>} y^{<2>} y^{<3>}) = p(y^{<1>}) p(y^{<2>} | y^{<1>}) \cdot p(y^{<3>} | y^{<1>} y^{<2>})$

4. Sampling Novel Sequences 采样新序列



$$p(y^{<1>} \dots y^{<Tx-1>})$$



以上, 构建了一个字级RNN (用的广义).

② Character-Level Language Model 字符RNN

Vocabulary = [a, aaron, .., zulu, <UNK>]

Vocabulary = [a, b, c, .., Z, L, ., , ;, 0..9, A..Z]

优点: 类 "Man" 未知单词. 可用单个字符表示根序.

缺点: 序列更长. 计算更复杂 (不如字级RNN)

⇒ 该字符RNN用的较少. (特殊情况下使用: 许多未知词)

小结: 建立一个RNN来观察英文文本.

构建单词级别的和字符级别的模型.

在你培训的语言模型中采样.

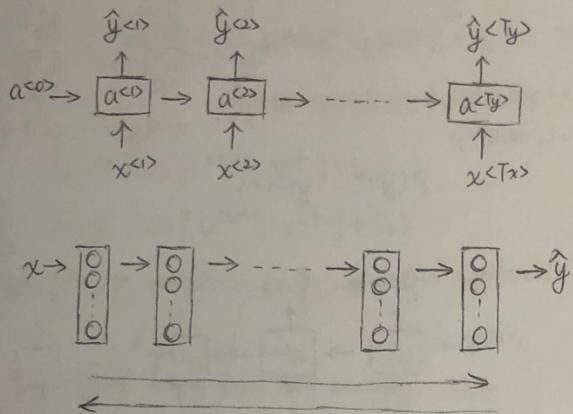
语言模型中的文本示例: News | Shakespeare.

(有文化级别的语言模型)

⇒ 如果该模型是在新闻文章中训练的, 它会生成左边.

⇒ 如果该模型在 Shakespeare 上训练的, 生成右~文本.

5. Vanishing Gradients with RNNs.



假设这是个很深的神经网络，从 \hat{y} 得梯度
很难反向传播去影响前期层的权重。

RNN：也存在类似问题。（很难由前面单词 \rightarrow 判断：单复数）

[The cat which ---, was full.
The cats which ---, were full.]

RNN模型：有附近效应。（Output \hat{y} 主要受接近 \hat{y} 值的影响）

这里的Output \hat{y} 很难被Input x 影响到。因为~~因为~~

Input x 是非常早期的序列。所以无论Output \hat{y} 是多少。

都不可通过反向传播到序列的开始，从而去修改NN的前期序列。

RNN：不擅长捕获远端依赖关系。

Nan：梯度爆炸。

Solution: Gradient Clipping

梯度向量如果大于某个临界值，要重新缩放某些向量，使其不那么大。

T_u 门控很容易设为0。只要“ $W_u \cdot \cdot + b_u$ ”是很大的负值。
在这附近的值更新门控值基本上为0。（非常接近0）。

$T_u = 0.000001$ 不要梯度弥散。因为 $T_u \approx 0$ 时， $C^{<t>} = C^{<t-1>}$

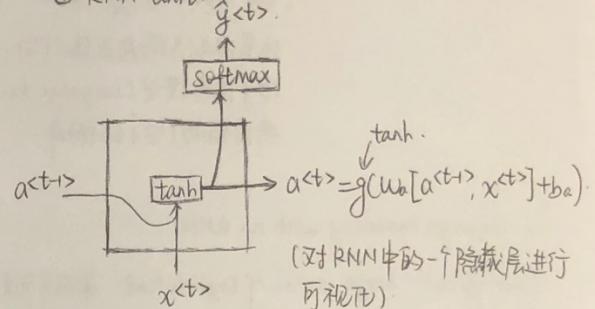
此时， C 保持很久不变。这可显著帮助梯度弥散~允许NN有很长范围的依赖性。

$C^{<t>}$ 100维 $\Rightarrow T_u$ 100维，近似取0.1

6. Gated Recurrent Unit (GRU) 门控循环单元。

GRU修改了RNN的隐藏层，从而更好地捕捉长距离关系，同时有助于减少梯度消失。

① RNN Unit.



② GRU

$u = \text{memory cell}$

$$C^{<t>} = a^{<t>}$$

$$\Rightarrow C^{<t>} = \tanh(W_o[C^{<t-1>}, x^{<t>}] + b_o)$$

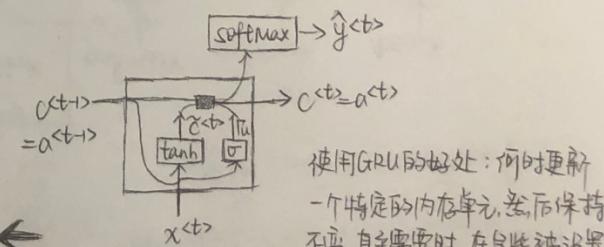
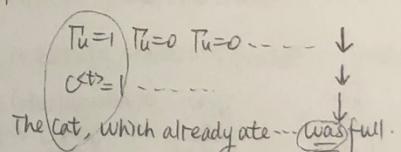
$$\Rightarrow T_u = \sigma(W_u[C^{<t-1>}, x^{<t>}] + b_u)$$

(Gamma Update)



$$C^{<t>} = T_u * C^{<t-1>} + (1 - T_u) * C^{<t-1>}$$

element-wise (元素乘法)



使用GRU的好处：何时更新一个特定的内存单元，之后保持不变。直至需要时，在某些被设置的值也可从记忆单元取出。

③ Full GRU.

$\tilde{C}^{t>} = \tanh(W_C [P_f * C^{t-1}, x^{t>}] + b_C)$. 计算记忆单元的候选新值.
 $\bar{C}^{t>}$ 相关性. 如何通过 C 来计算 $\bar{C}^{t>}$.

u. $P_u = \sigma(W_u [C^{t-1}, x^{t>}] + b_u)$.

f. $P_f = \sigma(W_f [C^{t-1}, x^{t>}] + b_f)$.

t. $C^{t>} = P_u * \tilde{C}^{t>} + (1 - P_u) * C^{t-1}$

7. LSTM (Long Short Term Memory) Unit.

① GRU and LSTM.

GRU (更简单, 可扩展性有利于构建较大的模型)

$$\tilde{C}^{t>} = \tanh(W_C [P_f * C^{t-1}, x^{t>}] + b_C)$$

$$\left\{ \begin{array}{l} P_u = \sigma(W_u [C^{t-1}, x^{t>}] + b_u) \\ P_f = \sigma(W_f [C^{t-1}, x^{t>}] + b_f) \end{array} \right.$$

$$C^{t>} = P_u * \tilde{C}^{t>} + (1 - P_u) * C^{t-1}$$

$$a^{t>} = \underbrace{C^{t>}}_{T_f}$$

LSTM (更大和灵活, 有三个门控)

$$\tilde{C}^{t>} = \tanh(W_C [a^{t-1}, x^{t>}] + b_C)$$

$$(update) \quad P_u = \sigma(W_u [a^{t-1}, x^{t>}] + b_u)$$

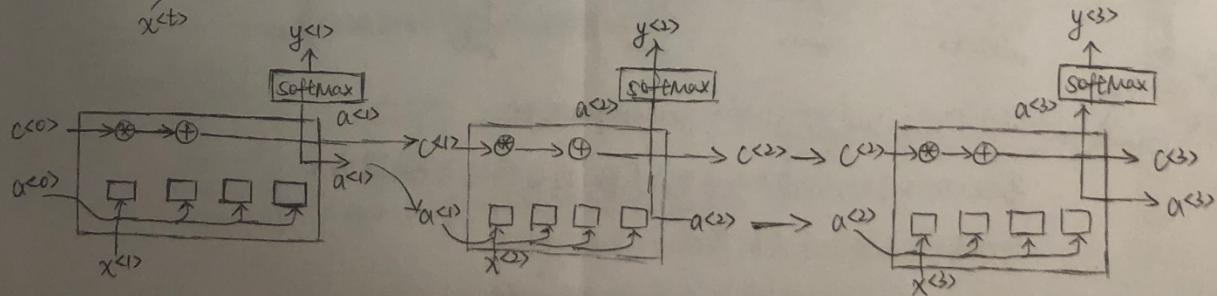
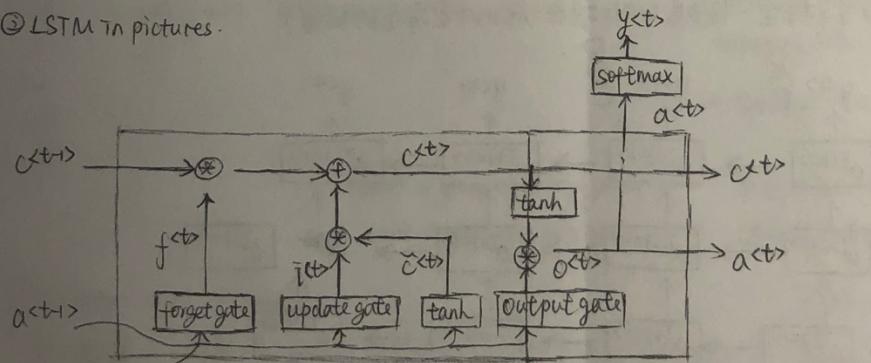
$$(forget) \quad P_f = \sigma(W_f [a^{t-1}, x^{t>}] + b_f)$$

$$(output) \quad P_o = \sigma(W_o [a^{t-1}, x^{t>}] + b_o)$$

$$C^{t>} = P_u * \tilde{C}^{t>} + P_f * C^{t-1}$$

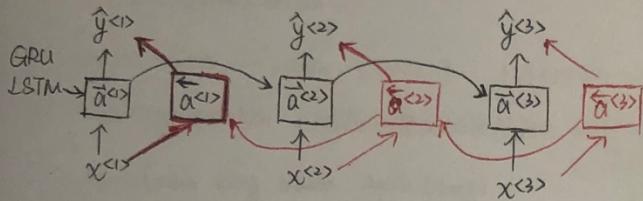
$$a^{t>} = P_o * \tanh(C^{t>})$$

② LSTM in pictures.



四. Bidirectional RNN 双向循环神经网络 (BRNN)

在一个时间点获得序列中前部分的和后部分的信息,

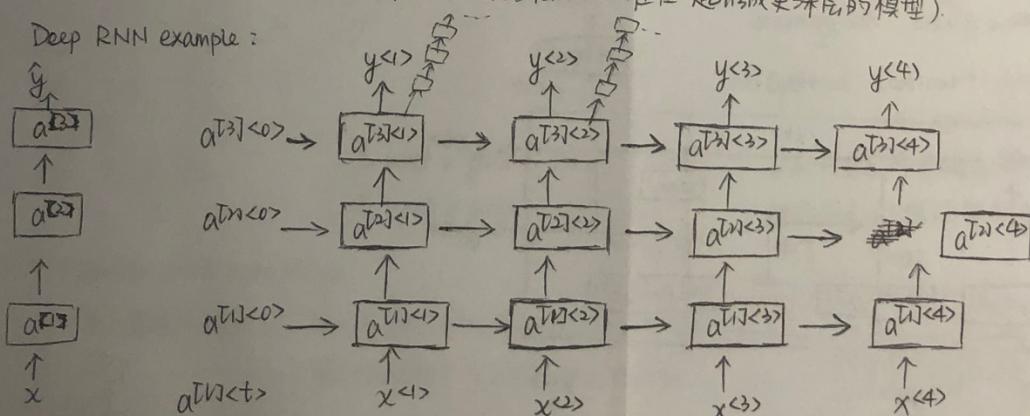


- 对于处理 NLP 问题的大段文本，带有 LSTM 的双向 RNN 是非常常用的。(可向前、可向后)。
- 对基本的 RNN 结构，或 GRU 或 LSTM 进行的修改，可以用来预测任何地方。在序列的中间可利用来自整个序列的信息。
- BRNN 缺点：需要整个序列，才可在任何地方预测。

例：建立语音识别系统。BRNN 会让你考虑整个演讲内容。
但如果是前向实现，需要等人停止说话，得到整个
讲话才可以实际处理它，并进行语音识别预测。

五. Deep RNNs. (非常复杂的学习模型：将多层 RNNs 堆在一起形成更深层的模型)

Deep RNN example:



$$a^{T21}{}^{<3>} = g(Wa^{T21}[a^{T21}{}^{<2>}, a^{T21}{}^{<3>}] + b_a^{T21})$$

当加入时间 t 时，即使只有很少的层数，这些网络也变得很大。

可以有一些无横向连接的模块。

Chapter 2. Word Representation

1. Introduction to Word Embeddings.

1. Word Representation.

$$V = [a, aaron, \dots, zulu, \text{UNK}] \quad |V| = 10000.$$

① 1-hot representation. (独热向量间的乘积为0. 无法度量任意两个向量之间的距离)

Man Woman King Queen Apple Orange
 (5371) (9853) (4914) (7157) (456) (625).

$$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

无法知晓

Apple 与 Orange 之间的相似度
是否高于 King, Queen, ...

② Featureized-representation = word embedding.

(特征化表示：文字嵌入)

	Man	Woman	King	Queen	Apple	Orange
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.17	0.69	0.03	-0.02
Food	0.09	0.01	0.02	0.01	0.95	0.97
---	---	---	---	---	---	---
Size	1	1	1	1	1	1
Cost	1	1	1	1	1	1

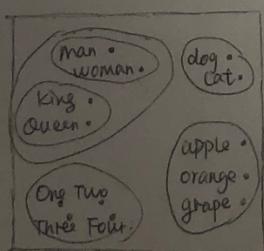
I want a glass of orange juice.

I want a glass of apple juice.

③ Visualize word embeddings.

t-SNE 算法(词嵌入): 可以学习相似的特征对于相关联的概念。

最终映射到一个更相似的特征向量。

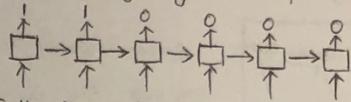


300 维特征空间
 Apple Orange (将它们嵌入在300维特征空间不同的点上)

映射到更低维的空间，使其可画在2维空间。

2. Using Word Embeddings

① Named entity recognition example.



Sally Johnson is an orange farmer.
 Robert Lin is an apple farmer.

当学习单词嵌入的算法，可检查大量文本~
(未标记)



采取该词嵌入，将其应用到命名中，在识别任务中训练集会变小

例：Orange, Apple, Grape \Rightarrow Fruit.

听作迁移学习，将知识迁移至任务上。

如：命名实体识别。

对于这个任务，也许有较小的已标记的训练集。(双向 RNN)

② Transfer Learning and Word Embeddings

\Rightarrow Learn word embeddings from large text corpus.

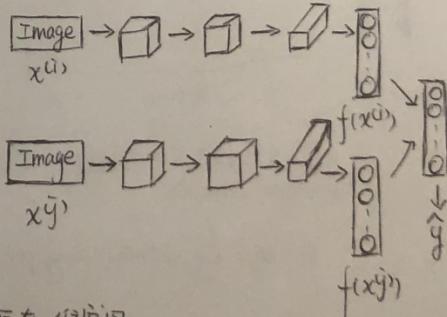
(Or download pre-trained embedding online.)

\Rightarrow Transfer embedding to new task with smaller training set. (100k words).

\Rightarrow Optional: Continue to finetune the word embeddings with new data.

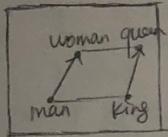
在新任务上训练模型时，可以选择性微调参数。
继续用新数据调整单词嵌入。

③ Relation to face coding (embedding)



3. Properties of Word Embeddings.

① Analogies using word vectors =



$$e_{\text{man}} - e_{\text{woman}} \approx e_{\text{king}} - e_{\text{queen}}$$

$$\arg \max_w \text{sim}(e_w, e_{\text{king}} - e_{\text{man}} + e_{\text{woman}})$$

Find word w.

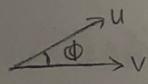
测量嵌入词 w 和等式右边量有多相似。

然后找到能最大化其相似度的词。

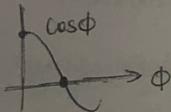
② Cosine Similarity

$$\text{sim}(e_w, e_{\text{king}} - e_{\text{man}} + e_{\text{woman}})$$

$$\text{sim}(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2} \quad (\text{U}, \text{V} \text{ 在产物}) \cdot (\text{U}, \text{V} \text{ 高度相似, 内在产物大}).$$

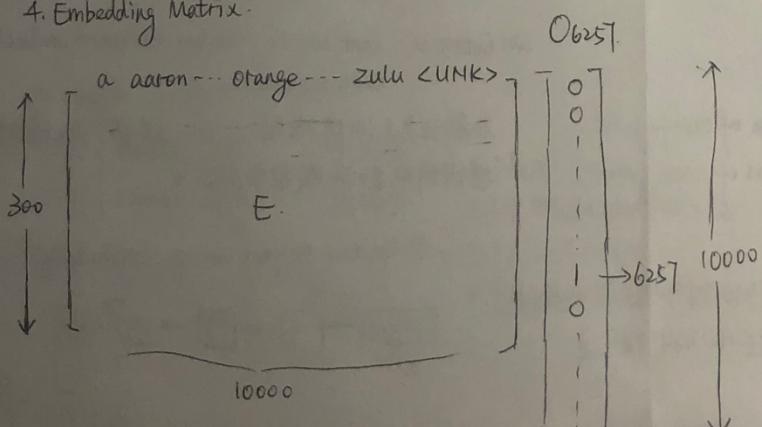


$\|u - v\|^2$. 欧几里得距离
是对非相似度的测量。



嵌入词类比相似概念: Man: Woman as Boy: Girl.
(类比推理系统)

4. Embedding Matrix.



$$E \cdot O_{6257} = e_{6257}$$

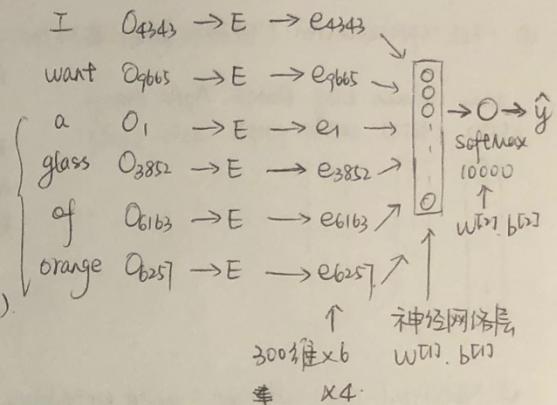
$$(300, 10K) \times (10K, 1)$$

$$E \cdot O_j = e_j \quad (\text{embedding for word } j)$$

2. Learning word embeddings.

1. Neural Language Model

I want a glass of orange juice.
4343 9665 1 3852 6163 6257



2. Other Context/ Target Pairs.

I want a glass of orange juice to go along with my cereal.
Target: cereal.

Context: Last 4 words.

构建语言模型: 用目标词前的几个词作上下文。

学习语言模型本身: 4 words on left & right.
(学习词嵌入)
Last 1 word.
Nearby 1 word.

3. Word2Vec

① Skip-grams.

I want a glass of orange juice to go along with my cereal.

Context Target

建立监督式学习. (输入单词, 会试图跳过一些单词
去预测另一些单词)

② Model:

Vocabulary Size = 10000k. $x \rightarrow y$

Context. c ("orange") \rightarrow Target t ("juice")
6257 4834

$c_c \rightarrow E \rightarrow e_c \rightarrow \text{SoftMax} \rightarrow \hat{y}$
($e_c = E_{c_0}$)

SoftMax: $p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10000} e^{\theta_j^T e_c}}$ \Rightarrow 计算缓慢.
(解决方法: 分层 SoftMax.) \Rightarrow 树上的每个内部节点, 都可用一个二元分类器表示.

$L(\hat{y}, y) = -\sum_{i=1}^{10000} y_i \log \hat{y}_i$, θ_t = parameter associate with output.

$y = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \leftarrow 4834$.
优化损失函数的话, 可得
到一系列相当好的嵌入
矢量.

不完美平衡树.

对于常见词: 只需几次遍历.

对于陌生词: 隐蔽深
(偏僻).

③ How to Sample the context C?

{ the, of, and ... 出现多 \Rightarrow 并不是全从训练
orange, apple, ... 出现少. 文本集中随机均匀采样的.
 \Rightarrow 有不同的启发式可用
来平衡常见词汇, 罕见词汇.