

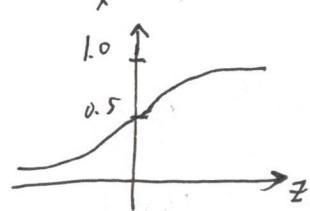
Deep learning

WEEK Course 1: Neural network and deep learning

Week 2: Basics of NN programming.

$$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix}, Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}] \quad |_{X \in M}$$

$n_x \times M$



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

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

~~cost func.~~ $\hat{y}^{(i)} = \sigma(w^T x^{(i)} + b), \sigma(z) = \frac{1}{1 + e^{-z}}$

Loss func: $L(\hat{y}, y) = -y \log(\hat{y}) - (1-y) \log(1-\hat{y})$

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

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

Gradient Descent:

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

computation graph (参考 2.7).

2.9 Logistic regression Gradient Descent

$$Z = w_1 x_1 + w_2 x_2 + b$$

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

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

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

该单个梯度，其 cost func:

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

a 为真实值， y 为预测值

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

$$\frac{\partial L(a, y)}{\partial a} = -\frac{y}{a} + \frac{1-y}{1-a}$$

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z}, \text{ 且 } \frac{\partial a}{\partial z} = a \cdot (1-a) \quad \downarrow a = \sigma(z)$$

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} = \left(-\frac{y}{a} + \frac{1-y}{1-a}\right) \cdot a(1-a) = a - y$$

$$\frac{\partial L}{\partial w_1} = x_1 \cdot dz$$

$$\frac{\partial L}{\partial w_2} = x_2 \cdot dz$$

$$\frac{\partial L}{\partial b} = dz$$

使用 $dw_1 = x_1 \cdot dz$, $dw_2 = x_2 \cdot dz$, $db = dz$

$$\text{Eqn } w_1 = w_1 - \alpha dw_1$$

$$w_2 = w_2 - \alpha dw_2$$

$$b = b - \alpha db$$

2.11 vectorization.

2.18. Explanation of logistic regression cost fun.

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

$$\sigma(z) = \sigma(w^T x + b) = \frac{1}{1 + e^{-z}}$$

$\hat{y} = P(Y=1|X)$, 给定样本 X 的情况下 $Y=1$ 的概率

if $y=1$: $P(y|x) = \hat{y}$

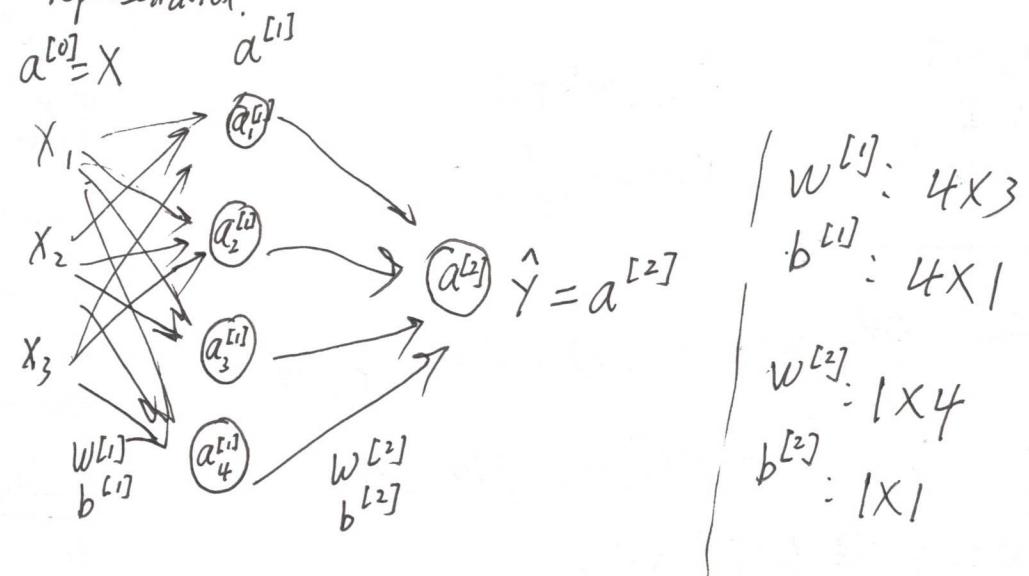
if $y=0$: $P(y|x) = 1 - \hat{y}$

令 y 上面的 \hat{y} 表示条件概率公式: $P(y|x) = \hat{y}^y (1-\hat{y})^{(1-y)}$
 \log 是单调增加的, $\log(\hat{y}^y (1-\hat{y})^{(1-y)}) \Rightarrow$
 $y \log \hat{y} + (1-y) \log(1-\hat{y})$

week 3. shallow neural networks

3.1 NN overview.

3.2 NN representation.



3.3 Computing a NN's output

$$\text{if } a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{bmatrix}, \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix} = \underbrace{\begin{bmatrix} \dots & w_{1,1}^{[1]} & \dots \\ \dots & w_{2,1}^{[1]} & \dots \\ \dots & w_{3,1}^{[1]} & \dots \\ \dots & w_{4,1}^{[1]} & \dots \end{bmatrix}}_{w^{[1]}} * \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix}$$

$$z^{[1]} = w^{[1]}x + b^{[1]}, \quad a^{[1]} = \sigma(z^{[1]}), \quad z^{[2]} = w^{[2]}z^{[1]} + b^{[2]}, \quad a^{[2]} = \sigma(z^{[2]})$$

3.4 Vectorizing across multiple examples.

若有m个样本，

用第一个样本为 $x^{[1]}$ 计算预测值 $\hat{y}^{[1]}$

$$= \dots x^{[2]} \quad \hat{y}^{[2]}$$

$$\text{第 } m \quad x^{[m]} \quad \hat{y}^{[m]}$$

$$a^{[2](1)}, a^{2}, \dots$$

$a^{[2](i)}$, (i)是第i个训练样本, [2]是第二层

向量化：

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

$$Z^{[1]} = \begin{bmatrix} z^{1} & z^{[1](2)} & \dots & z^{[1](m)} \\ \vdots & \vdots & \ddots & \vdots \end{bmatrix}$$

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

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

3.6 Activation functions.

之前已经用过 sigmoid 和 ReLU 作为 activation func.

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

可使用不同的 $g(z)$, 如 tanh, 双曲正切函数

$a = \tanh(z)$ 的值域 $+(-1)$

$$a = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}, \text{ 值域: } +(-1)$$

$g(z)$ = $\tanh(z)$ 优于 sigmoid, 因其梯度在值更接近0和0.5
三分之二处远比 sigmoid 平缓。

sigmoid, tanh 的共同缺点：z 特大 or 特小时，梯度小(≈ 0), ↓ gradient descent 的速度。

ReLU 函数： $a = \max(0, z)$

选择 activation function 理由：

$z \neq 0$ or 0

- ① 若输出是二分类问题，则输出用 sigmoid，其他单元用 ReLU。
- ② 在不精确时，用 ReLU，有时也用 tanh，但 ReLU 优点是 z 为负时，导数为0。
- ③ 另一缺点是 ReLU, Leaky ReLU, z 为负时导数为0。

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

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

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

$$\text{sigmoid, } a = \frac{1}{1+e^{-z}}$$

$$\text{tanh, } a = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\text{ReLU, } a = \max(0, z)$$

$$\text{leaky ReLU, } a = \max(0, \alpha z)$$

结论：

① Relu, Leaky Relu 的导数都 > 0 , 只需 if-else 语句, sigmoid 导数立运算, 前者更快.

② sigmoid, tanh 在正负饱和区的梯度 ≈ 0 , 造成梯度弥散, Relu, 和 Leaky Relu 已经大于 0 部分导数为常数, 不会弥散, Relu 在负半区梯度为 0, 神经元此时不训练, 产生稀疏性, Leaky Relu 不会.

z 在 Relu 的负一半为 0, θ 在 hidden layer 足够大使得 $z > 0$, 很快.

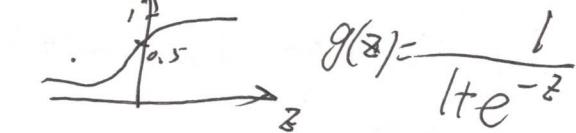
3.7. Why need a nonlinear activation function.

反面言之：若选择 $f(z)$ 在 hidden layer 用线性 activation func, 在输出层用 sigmoid func, 那么这个模型的泛化能力和没有使用 hidden layer 相当 logistic regression 一样.

活性函数的组合还是线性函数, 比如 $3/4 \sin(4\pi z)$, 否则无法训练有意义的模型.

3.8 Derivatives of activation func: a

1) sigmoid activation func:



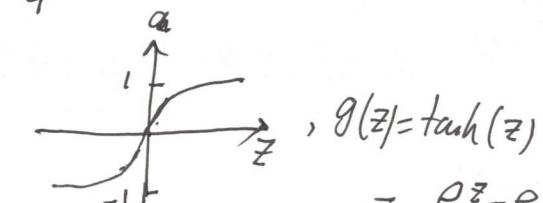
$$g(z) = \frac{1}{1 + e^{-z}}$$

$$\frac{d g(z)}{dz} = -\frac{e^{-z}}{(1 + e^{-z})^2} = \frac{1}{1 + e^{-z}} \left(1 - \frac{1}{1 + e^{-z}}\right) = g(z)(1 - g(z))$$

$z = \pm 10$ 时, $\frac{d g(z)}{dz} \approx 0$

$$z = 0 \text{ 时}, \frac{d g(z)}{dz} = \frac{1}{4}$$

2) Tanh activation func.



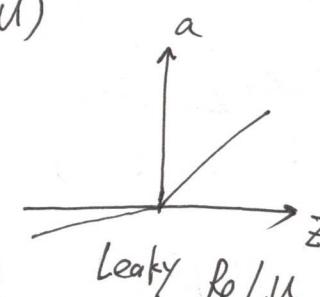
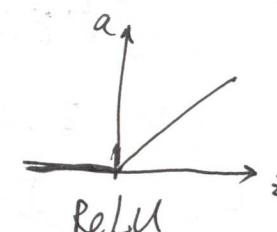
$$g(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

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

$z = \pm 10$, $\frac{d g(z)}{dz} \approx 0$

$$z = 0, \frac{d g(z)}{dz} = 1 - 0 = 1$$

3) Rectified Linear Unit (ReLU)



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

$$g(z)' = \begin{cases} 0 & , \text{if } z < 0 \\ 1 & , \text{if } z > 0 \\ \text{undefined} & , \text{if } z = 0 \end{cases}$$

$z = 0$ 时, z 不等于 0 或 1

④ 2.9.

4) Leaky linear unit (Leaky ReLU)

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

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

$z=0$ At, $g(z) \neq 0.01$ or !

back propagation:

$$\textcircled{1} dz^{[2]} = A^{[2]} - Y, Y = [Y^{[1]}, Y^{[2]}, \dots, Y^{[m]}]$$

$$\textcircled{2} dW^{[2]} = \frac{1}{m} dz^{[2]} A^{[1]T}$$

$$\textcircled{3} db^{[2]} = \frac{1}{m} \text{np.sum}(dz^{[2]}, \text{axis}=1, \text{keepdims=True}) \quad \text{189}$$

$$\textcircled{4} dz^{[1]} = \underbrace{W^{[2]T} dz^{[2]}}_{(n^{[1]}, m)} * \cdot g^{[1]'} * \underbrace{(z^{[1]})}_{\text{activation func. of hidden layer}} \quad (n^{[1]}, m)$$

$$\textcircled{5} dW^{[1]} = \frac{1}{m} dz^{[1]} x^T$$

$$\textcircled{6} db^{[1]} = \frac{1}{m} \text{np.sum}(dz^{[1]}, \text{axis}=1, \text{keepdims=True}) \quad (n^{[1]}, 1)$$

$$\leftarrow dz = a - y$$

$$\leftarrow dw_2 = x_2 \cdot dz$$

$$\leftarrow db = dz$$

$$\leftarrow dw_1 = x_1 \cdot dz$$

189 推导

3.10 Backpropagation Intuition.

$$\begin{matrix} x \\ w \\ b \end{matrix} \Rightarrow z = w^T x + b \Rightarrow a = g(z) \Rightarrow L(a, y)$$

$$\frac{da}{d\alpha} = \frac{dL(a, y)}{d\alpha} = (-y \log \alpha - (1-y) \log(1-\alpha))' = -\frac{y}{\alpha} + \frac{1-y}{1-\alpha}$$

$$dz^{[2]} = a^{[2]} - y, dW^{[2]} = dz^{[2]} a^{[1]T}, db^{[2]} = dz^{[2]}$$

$$dz^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]}) .$$

$$dW^{[1]} = dz^{[1]} x^T, db^{[1]} = dz^{[1]}$$

forward propagation:

$$\textcircled{1} z^{[1]} = W^{[1]} x + b^{[1]}$$

$$\textcircled{2} a^{[1]} = g(z^{[1]})$$

$$\textcircled{3} z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$

$$\textcircled{4} a^{[2]} = g^{[2]}(z^{[2]}) = \sigma(z^{[2]})$$

3.11 Random Initialization

若 w 全初始化为 0，则各层单元皆同。

随机初始化：

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

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

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

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

为什么 $* 0.01$ ，因为若 w 很大， z 很大， \tanh or sigmoid 在平滑部，梯度下降慢。

Week 4 Deep Neural Networks

4.1 Deep L-layer neural network

符号之：L：层数

输入层不包含在 L 中，为 0 层

$n^{[L]}$ 为第 L 层层数

4.2 Forward and backward propagation.

forward propagation

$$z^{[l]} = W^{[l]}.a^{[l-1]} + b^{[l]}, \quad a^{[l]} = g^{[l]}(z^{[l]})$$

$$\text{向量化: } z^{[l]} = W^{[l]}.A^{[l-1]} + b^{[l]}, \quad A^{[l]} = g^{[l]}(z^{[l]})$$

back propagation:

? 待补

$$\begin{cases} \textcircled{1} \quad dz^{[l]} = da^{[l]} \cdot g'^{[l]}(z^{[l]}) \\ \textcircled{2} \quad dw^{[l]} = dz^{[l]}.a^{[l-1]} \\ \textcircled{3} \quad db^{[l]} = dz^{[l]} \\ \textcircled{4} \quad da^{[l-1]} = W^{[l]T}.dz^{[l]} \\ \textcircled{5} \quad dz^{[l]} = W^{[l+1]T}da^{[l+1]}.g'^{[l]}(z^{[l]}) \\ (\textcircled{5} \text{ 由 } \textcircled{4} \text{ 及 } \textcircled{1} \text{ 推导}) \end{cases}$$

$$\text{向量化: } \textcircled{1} \quad dz^{[l]} = dA^{[l]} \cdot g'^{[l]}(z^{[l]})$$

$$\textcircled{2} \quad dw^{[l]} = \frac{1}{m} dz^{[l]}.A^{[l-1]T}$$

$$\textcircled{3} \quad db^{[l]} = \frac{1}{m} \text{np.sum}(dz^{[l]}, axis=1, keepdims=True)$$

$$\textcircled{4} \quad dA^{[l-1]} = W^{[l]T}.dz^{[l]}$$

4.3 Forward propagation in a Deep network.

4.4 Getting your matrix dimensions right

结论：

$$w^{[l]}: (n^{[l]}, n^{[l+1]})$$

$$b^{[l]}: (n^{[l]}, 1)$$

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

$$dw^{[l]}, w^{[l]} \text{ 为矩阵}$$

$$db^{[l]}, b^{[l]} \text{ - 扑月。}$$

w, b 为矢量，但 z, a 为 $n \times m$ 矩阵。

从右往左向量化后修改。

右括号， $m + 1$ 拼串。

$$z^{[l]} = (z^{[l][1]}, z^{[l][2]}, \dots, z^{[l][m]}) \\ (n^{[l]}, m)$$

$$A^{[l]}: (n^{[l]}, m), A^{[0]} = X = (n^{[0]}, m)$$

$$z^{[l]} = w^{[l]} \cdot X + b^{[l]}$$

4.5 Why Deep representation?

前面几层学简单加，之后的层学更复杂的特征。
若不加层级，就要搞极低 \uparrow neuron 数量

4.6 Building blocks of deep neural networks.

4.7 Parameters vs Hyperparameters

Hyperparameters: learning rate α , Iteration, L (Number of hidden layers)
 $n^{[l]}$ (Number of neurons in one hidden layer), choice of activation function.
 内部参数：idea - code - experiment - idea 试验

第二门课：Improving Deep NN: Hyperparameter tuning, Regularization and optimization.

Week 1: practical aspects of deep learning

train set, dev set, test set . 60%, 20%, 20%
(cv set)

两者相对大，test set 可小一些

1.2 Bias / Variance., 看 notes

如何同时 High Bias, high variance.

1.3 Basic Recipe for Machine Learning

训练模型时，先拉直 bias，把 bias ↓ 则能
能拟合，再想法（更复杂模型，之后反向梯度）
使 Var ↓ Variance.

1.4 Regularization.

解决 overfitting (high variance) ① regularization
② 正则化法。

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

$$\text{L}_2 \text{ regulation: } \|w\|_2^2 = \sum_{j=1}^n w_j^2 = w^T w \leftarrow \text{frobenius norm}$$

$$\text{L}_1 \text{ reg - : } \|w\|_1$$

$$\|\cdot\|_2^2 \text{ 为 Frobenius norm}$$

$$\begin{aligned} \text{GD: } w^{[t+1]} &= w^{[t]} - \alpha \left[\text{from back propa.} + \frac{\lambda}{n} w^{[t]} \right] \\ &= \left(1 - \frac{\alpha \lambda}{n}\right) w^{[t]} - \alpha \cdot \text{from back propa.} \end{aligned}$$

Regularized logistic regression: (参考 ML 29.7.3).

$$\theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right], \text{ for } j \neq 0.$$

1.5 why regularization reduces overfitting

入 ↑, w ↓, z ↑, z 值越大

使用网络简单。(or, 使用层的数目 ↓), tanh 在某些情况下，
其值 ↓ Variance.

1.6 Drop out regularization.

x_1	0	0	.
x_2	0	0	0
x_3	0	0	.
x_4	0	0	.

在训练 NN overfitting, 因 dropout 方法，
遍历每一层，设置消除部分神经元使之不被选中
从而消除一些节点，从而降低整个的 NN
获得概率分布 d_3 后，
 $a_3 = d_3 \cdot a_3$
使 a_3 中概率变均匀。

$$\text{最后: } a_3 = \frac{\alpha_3}{\text{keep-prob}}$$

使 Z^L 的期望不变即 1.0.

训练阶段不用 drop out 方案.

1.7 understanding Dropout

$$\begin{array}{cccccc}
 & 0 & 0 & & w^{[1]}: 7 \times 3 \\
 x_1 & 0 & 0 & & w^{[2]}: 7 \times 7 \\
 x_2 & 0 & 0 & 0 & w^{[3]}: 3 \times 7 \\
 x_3 & 0 & 0 & 0 & 0 \\
 & 0 & 0 & 0 & 0 \\
 & 0 & 0 & 0 & 0 \\
 & 0 & 0 & 0 & 0
 \end{array}$$

$w^{[2]}$ 权重过大, (7×7) , 为

overfitting, keep-prob 低一些, 0.5.

若其他层 overfitting 不严重, keep-prob 高一些, 0.7.

若不 overfitting, —— 1.

1.8 Other regularization methods

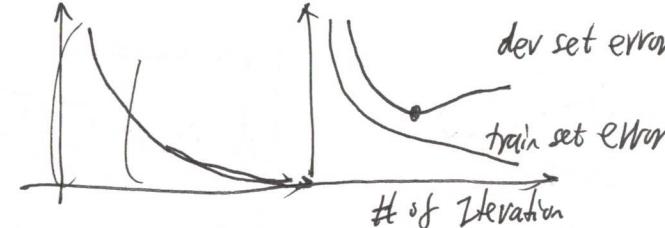
(1). Data augmentation.

图片翻转, 增加
扭曲.

增加背景图片, 很难发现

(2) early stopping

训练时先固定 w 的值, 但随着训练 T , w 1



early stopping 就是在 w 中
大于时停止.

“正则化”：这是一个时期 D. 为一个任务。如防 overfitting，只关注 w, b ，使 $J(w, b)$ 变小；防 high variance，再用另一套工具实现。

early stopping 的缺点：不能解决处理上两个问题。

优点：D. 运行一次 Gradient Descent，输出 w 较小值。

L2 regularization 缺点：计算时间长，尝试很多种参数。

1.9 Normalizing Inputs

归一化输入可以加速梯度下降 model.

假设只有两个特征，normalization 两步

① subtract out or to zero out the mean. $\bar{x} = \frac{1}{m} \sum_{i=1}^m x^{(i)}$

② normalize variance : $\sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \bar{x})^2$ $x := x - \bar{x}$

$$x = \frac{x - \bar{x}}{\sigma^2}$$

在 train set 上得到 \bar{x}, σ^2 也应用到 dev set 上。

1.10 Vanishing / Exploding gradients.

导致梯度衰减有时权重的值非常大，也非常小。

$$\begin{array}{cccc} x_1 & 0 & 0 & 0 \\ x_2 & 0 & 0 & 0 \end{array} \quad \cdots \quad \begin{array}{c} 0 \rightarrow 0 \rightarrow \hat{y} \\ 0 \rightarrow 0 \rightarrow \hat{y} \\ w^{[1]} \quad w^{[2]} \end{array}$$

$\text{fix } g(z) = z, b=0$

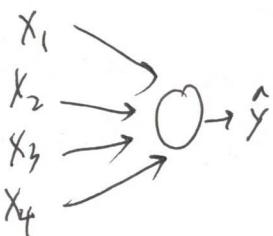
$$\hat{y} = w^{[L]} w^{[L-1]} \cdots w^{[3]} w^{[2]} w^{[1]} x$$

If $w^{[i]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$, then $\hat{y} = w^{[L]} \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{L-1} x$ explode

If $w^{[i]} = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$, $\hat{y} = w^{[L]} \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}^{L-1} x$ Vanish

1.11 Weight Initialization for deep networks.

- # neurons to init



$$z = w_1 x_1 + w_2 x_2 + \cdots + w_n x_n, b=0$$

为了避免 z 过大 or 过小, 设 $w_i = \frac{1}{n}$
 n 为输入特征数量. $\text{Var}(w) = \frac{1}{n}$

$w^{[L]} = np.random.rand(shape) * np.sqrt(\frac{1}{n^{L-1}})$

If ReLU, $g^{[L]}(z) = \text{ReLU}(z)$, $\text{np.sqrt}(\frac{2}{n^{L-1}})$

$$\text{variance} = \frac{2}{n}$$

~~to fanin $\sqrt{\frac{1}{n^{L-1}}}$ 转接 np.sqrt(2/n)~~

~~to tanh actv, $\sqrt{\frac{1}{n^{L-1}}}$ 转接 np.sqrt(1/n), Xavier init~~

Yoshua Bengio: $\sqrt{\frac{2}{n^{L-1} + n^L}}$

1.12 Numerical approximation of gradients

$$\frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon}$$

双边差比单边差, 确认 $g(\theta)$ 是不是 $f(\theta)$
导数的正确实现.

1.13 Gradient checking

$$w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]} \rightarrow \text{grad}, \theta$$

$$dw^{[1]}, db^{[1]}, \dots, dw^{[L]}, db^{[L]} \rightarrow \text{grad}, d\theta$$

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

计算 $\frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\| + \|d\theta\|_2} \leftarrow$ 没差平分和
 \leftarrow 平方根

$$\epsilon = 10^{-7}, \quad \downarrow \approx 10^{-7} \rightarrow \text{good}$$

$$\approx 10^{-5} \rightarrow \text{挂壁}$$

$$\approx 10^{-3} \rightarrow \text{worry}$$

1.14 Gradient checking implementation notes.

- ① 在训练中用 gradient checking, 不用梯度
- ② 若 gradient checking 失效, 考虑梯度的范数
- ③ 注意正则化
- ④ Gradient checking & dropout 不会同时用
- ⑤ ...

mini-batch $X^{t,y}, Y^{t,y}$
repeat
For $t = 1, 2, \dots, 5000$
{ forward prop. on $X^{t,y}$
 $Z^{[l]} = W^{[l]} X^{t,y} + b^{[l]}$
 $A^{[l]} = g^{[l]}(Z^{[l]})$
 $\hat{A}^{[L]} = g^{[L]}(Z^{[L]})$
 $\frac{\lambda}{2 \cdot 1000} \sum \|W^{[l]}\|_F^2$
compute cost. $J = \frac{1}{1000} \sum_{i=1}^t L(\hat{Y}^{(i)}, Y^{(i)}) +$
Backprop to compute gradients $J^{t,y}$ (using $(X^{t,y}, Y^{t,y})$)
 $w^{[L]} := w^{[L]} - \alpha dW^{[L]}$
 $b^{[L]} := b^{[L]} - \alpha db^{[L]}$

Week 2 Optimization algorithms

2.1 Mini-batch Gradient descent

$$X = [x^{(1)} \ x^{(2)} \ \dots \ x^{(m)}], (n_x, m)$$

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

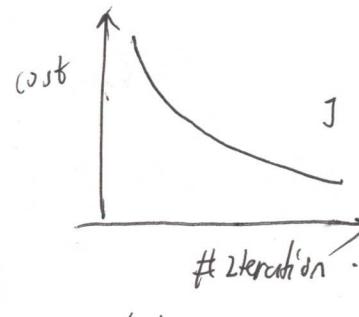
把 n_x 分成 $\frac{n_x}{2,000}$ 份，每份为 mini-batch

假设 500 例样本，每个 mini-batch 为 1000，共有 500 个 mini-batch

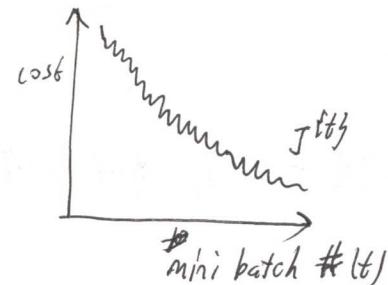
$$X = [x^{(1)} \ x^{(2)} \ x^{(3)} \ \dots \ x^{(1000)} \ x^{(1001)} \ \dots \ x^{(2000)} \ \dots \ x^{(5000)}]$$

$$Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(1000)} \ y^{(1001)} \ \dots \ y^{(2000)} \ \dots \ y^{(5000)}]$$

2.2 Understanding Mini-batch gradient descent



Batch AD



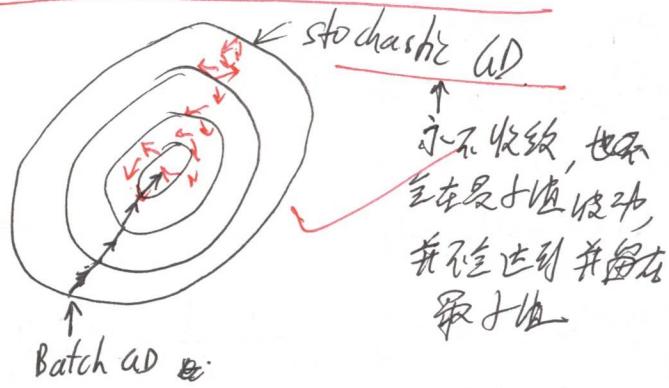
Mini-batch AD

"1 epoch" \leftarrow single pass through training set.

应该用 mini-batch AD, (not batch)

若行训练样本数为 n, BP 为 Batch AD.

若 n=1, 则为 stochastic AD.



mini-Batch 大约在 $1 \sim n$ 之间。

若 mini-batch AD 的范围很窄, ↓ learning rate

若样本大于 2000, 直接用 Batch AD

样本大, minibatch 在 64~512 之间
($\times 128, 256, 512$)

2.3 Exponentially weighted averages

温度的例子

V_0

$V_t = 0.9V_{t-1} + 0.1\theta$, 前一天 0.9, 今天 0.1

:

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

$\uparrow \beta \quad \uparrow 1-\beta$

$$V_t = \beta V_{t-1} + (1-\beta)\theta_t \quad \text{可视为 } \frac{1}{1-\beta} \text{ 天的温度}$$

$\beta = 0.9, \frac{1}{1-\beta} = 10$, 也就说 +10 天的平均值

$\beta = 0.98, \frac{1}{1-\beta} = 50 \rightarrow -50 \text{ 天} -$

2.4 Understanding exponentially weighted averages. → 古内存

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

$$\theta_{100} = 0.9V_{99} + 0.1\theta_{100}$$

$$\theta_{99} = 0.9V_{98} + 0.1\theta_{99}$$

$$\theta_{98} = 0.9V_{97} + 0.1\theta_{98}$$

: 不迭代入

$$V_{100} = 0.1\theta_{100} + 0.1 \times 0.9\theta_{99} + 0.1 \times 0.9^2\theta_{98} + 0.1 \times 0.9^3\theta_{97} + 0.1 \times (0.9)^4\theta_{96} + \dots$$

exponential decay func: $0.1 \downarrow \dots$

相当于当 t 时的 V_t 与 exponential decay func 的乘积, 最后才加 θ_{100} , 得到 V_{100}

偏差修正: $0.1 \times 0.1 \times (0.9)^1 \times 0.1 \times (0.9)^2 \times 0.1 \times (0.9)^3 \times \dots$

代码中令 $V_0 = 0$ → 退化 /

$$V_0 := \beta v + (1-\beta)\theta_1$$

$$V_\theta := \beta v + (1-\beta)\theta_2$$

2.5 Bias correction in exponentially weighted averages.

$$\hat{V}_0 = 0 \text{ 时}, V_1 = 0.98 V_0 + 0.02 \theta = 0.020,$$

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

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

前面几天估计不准，偏大，

故用不带 V_0 ，改用 $\frac{V_t}{1-\beta^t}$

$$t=2 \text{ 时} \quad \frac{V_2}{1-0.98^2} = \frac{0.0196 \theta_1 + 0.02 \theta_2}{0.0396}$$

变成了偏差， θ_1

t 很大时， $\beta^t \approx 0$

2.6 Gradient Descent with momentum

(快于梯度下降)

其思想：计算梯度的指数滑动平均，并利用该梯度更新权重。

$$V_{dw} = \beta V_{dw} + (1-\beta) dw$$

$$V_{db} = \beta V_{db} + (1-\beta) db$$

$$W := W - \alpha V_{dw}$$

$$b := b - \alpha V_{db}$$



想成小波向流动
加快向右流动

回想之前无动量的例子，GD with momentum 算

是在估计全局梯度。

2.7 Root Mean Square prop. (RMSprop)

$$S_{dw} = \beta S_{dw} + (1-\beta) dw^2$$

$$S_{db} = \beta S_{db} + (1-\beta) db^2$$

$$w := w - \alpha \frac{dw}{\sqrt{S_{dw} + \epsilon}}$$

$$b := b - \alpha \frac{db}{\sqrt{S_{db} + \epsilon}}$$

优点：↓ 损耗，可使用较大的学习率 α ，加快算法

S_{dw} 及 S_{db} 取大， dw 及 db

S_{dw} 及 S_{db} 取小

w 及 b

为了避免 0，在 w, b 的更新中加入 $\epsilon \approx 10^{-8}$

2.8 Adam Optimization Algorithm

(Adaptive Momentum Estimation).

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

One iteration t :

compute dw, db using current mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1-\beta_1) dw, V_{db} = \beta_1 V_{db} + (1-\beta_1) db \leftarrow \text{momentum}$$

$$S_{dw} = \beta_2 S_{dw} + (1-\beta_2) dw^2, S_{db} = \beta_2 S_{db} + (1-\beta_2) db^2 \leftarrow \text{RMSprop}$$

$$\begin{cases} V_{dw}^{\text{corrected}} = \frac{V_{dw}}{(1-\beta_1^t)}, V_{db}^{\text{corrected}} = \frac{V_{db}}{(1-\beta_1^t)} \\ S_{dw}^{\text{corrected}} = \frac{S_{dw}}{(1-\beta_2^t)}, S_{db}^{\text{corrected}} = \frac{S_{db}}{(1-\beta_2^t)} \end{cases}$$

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

$$\text{其中 } \beta_1 = 0.9$$

$$\beta_2 = 0.999$$

Hyperparameters choice:

λ : needs to be fine

$$\beta_1: 0.9$$

$$\beta_2: 0.999$$

$$\epsilon: 10^{-8}$$

2.9 · Learning rate decay.

batch size 与 learning rate: $\text{batch size} \downarrow \rightarrow \text{learning rate}$.

gg is learning rate decay/decay

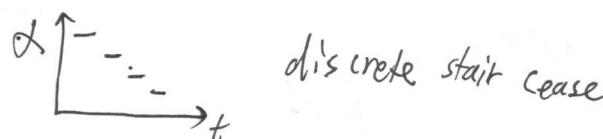
$$\text{新的 } \alpha = \frac{1}{\text{decay rate} * \text{epoch_num}} \alpha_0$$

\uparrow decay rate \uparrow mini-batch & no. of gg

other learning rate decay methods

$$\lambda = 0.95^{\text{epoch_num}} \cdot \alpha_0 \quad \text{exponential decay}$$

$$\alpha = \frac{k}{\sqrt{\text{epoch_num}}} \cdot \alpha_0 \quad \text{or} \quad \frac{k}{\sqrt{t}} \alpha_0$$



2.10 The problem of local optima.
→ no saddle point.

Week 3. Hyperparameter tuning, Batch normalization, and propagation framework

3.1 Tuning process

Hyperparameters:

$$\alpha, \beta_1, \beta_2, \epsilon$$

(0.9), (0.999), (10^{-8})

layers

hidden units

learning rate decay

mini-batch size.

选值: ① 通过乱序训练
② 使用反向传播梯度

3.2 Using an appropriate scale to pick hyperparameters.
约 0.001 ~ 0.1 适合梯度下降的支撑

3.3 Hyperparameters tuning in practice: Pandas vs. Caviar

→ 拾金，对 Q -> model 有影响

→ Q 对于 model 有影响同时执行

3.4 Normalizing activations in a network

既 normalize $\lambda x_1, x_2, \dots$
也 normalize $\beta z^{(1)}, z^{(2)}, \dots$

$$\mu = \frac{1}{m} \sum_i z^{(i)}$$

$$\sigma^2 = \frac{1}{m} \sum_i (z_i - \mu)^2$$

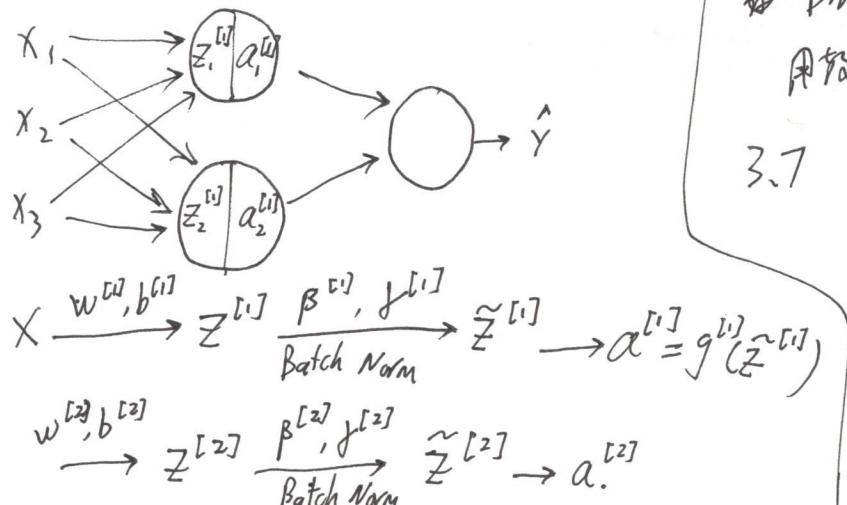
$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

这样 hidden units 都 $\sim N(0, 1)$, 但还要
其分布不那么相似,

$$\tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta, \quad \gamma, \beta \text{ 需训练}$$

$$\text{if } \gamma = \sqrt{\sigma^2 + \epsilon}, \beta = \mu, \text{ then } \tilde{z}^{(i)} = z^{(i)}$$

3.5 Fitting Batch norm into a NN.



in tensorflow, Batch norm: tf.nn.batch_normalization.

Batch Norm 要对每一个 mini-batch $x^{[i]}$ 处理

对每一个 mini-batch $x^{[i]}$, 计算 $\bar{z}^{[i]}$,

$$x^{[i]} \xrightarrow{w^{[i]}, b^{[i]}} z^{[i]} \xrightarrow{\gamma^{[i]}, \beta^{[i]}, \text{BN}} \tilde{z}^{[i]} \xrightarrow{g^{[i]}} g^{[i]}(\tilde{z}^{[i]}) = \alpha^{[i]} w^{[i]}, b^{[i]} \xrightarrow{\gamma^{[i]}, \beta^{[i]}} z^{[i]} \dots$$

这样每次 BN. $z^{[i]}$ 时, $(\gamma^{[i]} - \bar{z}^{[i]} = w^{[i]} a^{[i-1]} + b^{[i]})$, $b^{[i]}$ 都被
修正了, 可设其为 0, $z^{[i]} = w^{[i]} a^{[i-1]}$, $\tilde{z}^{[i]} = \gamma^{[i]} z^{[i]} + \beta^{[i]}$

3.6 Why does Batch Norm work?

① 对隐藏层归一化

② 减轻重比例的网络更薄后或更厚后,
“covariate shift”.

Batch BN 比 $\gamma z + \beta$ 好, 因为 $\gamma z + \beta$ 的值和方差不变, 往后层更坚挺.

BN 有正则化 regularization 效果, 因每次只对一个 mini-batch 算 μ, σ^2 有帮助.
用整个 mini-batch 会等效.

3.7 Batch Norm at test time.

BN 在 mini-batch 逐一批处理, 但测试时要对每个样例处理, 如何?

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

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

训练时， μ, σ^2 是对每个 mini-batch 的，
需要的 μ, σ^2 是根据训练集估计的。

通常用梯度下降平均来估计 μ, σ^2 。

3.8 Softmax regression

识别大于两个类，

识别别猫(0), 狗(1), 鸡(2), 其他(3)四类为例：

$X \rightarrow \dots$

$\hat{Y} = [y^{(1)} \ y^{(2)} \ y^{(3)} \dots y^{(m)}]$

$\hat{Y} = [P(\text{other}|x) \ P(\text{cat}|x) \ P(\text{dog}|x) \ P(\text{chicken}|x)]$

4个概率和为1。

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

Activation function:

$$t = \frac{e^{z^{[l]}}}{\sum_{j=1}^4 t_j}$$
$$a^{[l]} = \frac{e^{z^{[l]}}}{\sum_{j=1}^4 e^{z^{[l]}}}$$
$$z^{[l]} = \begin{bmatrix} 5 \\ 2 \\ 1 \\ 3 \end{bmatrix}$$
$$t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix}, \sum_{j=1}^4 t_j = 176.3$$
$$a^{[l]} = \frac{t}{176.3}$$

3.9 Training a softmax classifier.

hard max $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$
只在 $a^{[l]}$ 中有 1，其他为 0。

soft max: 柔软化

softmax 和 logistic regression 的区别 (C ≥ 2)

Loss func 能使输出更接近正确的类别，使该类别概率上升。

对于一个样本，loss func: $L(y, \hat{y}) = -\sum_{j=1}^C y_j \log \hat{y}_j$

整个训集 loss: $J(w^{[0]}, b^{[0]}, \dots) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$

$$Y = [y^{(1)} \ y^{(2)} \ y^{(3)} \dots y^{(m)}] \quad \hat{Y} = [\hat{y}^{(1)} \ \hat{y}^{(2)} \ \dots \hat{y}^{(m)}]$$
$$= \begin{bmatrix} 0 & 0 & 1 & \dots \\ 0 & 0 & 0 & \dots \\ 0 & 1 & 0 & \dots \\ 0 & 0 & 0 & \dots \end{bmatrix} \quad = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix} \quad \dots \quad (4, m) \quad (4, m)$$

反向传播一步: $d_z^{[l]} = \frac{\partial J}{\partial z^{[l]}} = \hat{y} - y$, 其他部分略去

3.10 Deep learning framework,

3.11 TensorFlow & Notes.

Course 3 Structuring Machine Learning projects.

流程：用 training set 训练不同 model，用 dev set 选模型，用 test set 评估 model.

dev set, 和 test set 都是同一分布， \leftarrow 把他们 shuffle.

1.6 Size of dev and test sets.

1.7 When to change dev/test sets and metrics.

当你的评估指标无法正确衡量算法之好坏时，这改变评估指标。

常用的错误率指标可以表示：Error = $\frac{1}{m_{dev}} \sum_{i=1}^{m_{dev}} I\{Y_{pred}^{(i)} \neq Y^{(i)}\}$

如上所示：

$$Error = \frac{1}{m_{dev}} \sum_{i=1}^{m_{dev}} w^{(i)} \cdot I\{X_{pred}^{(i)} \neq X^{(i)}\}$$

$X^{(i)}$ 是图片且 $w^{(i)} = 1$ ，是 pornographic 且 $w^{(i)} = 10$ or 100 ，惩罚大。

1.8 Why human-level performance?

Bayes optimal error：理论上可能达到的最小错误率
(有时写作 Bayesian)

1.9 Avoidable bias :

即 Bayes optimal error 与训练错误率之间的差值。

1.10 Understanding human-level performance.

理解 Bayes optimal error 的差距，意味着 avoidable bias 大

训练错误率和开发错误率的差距，差距越大，越容易避免。

Week 1 ML

1.1 Why ML strategy?

教书哪些改进策略更值得尝试

1.2 Orthogonalization. (正交化)
正交指：一个维度只和一个变量

1.3 Single # Evaluation metric

有一个单一指标评价 model 很重要，
precision

Recall

F₁ score

etc.

!:

1.4 Satisficing and optimizing metrics.

- 满足指标

其他作为满意度指标

1.5 Train / Dev / test distributions,

development set 也叫 hold out · cross validation set.

1.11 surpassing human-level performance.

1.12 Improving your model performance. Notes.

注意：

avoidable bias:

Week 2 ML Strategy 2.

2.1 Carrying out error analysis

Error analysis: 将错误划归到人类的表现，检查其中的错误.

看那些错的很离谱，分析那些错误大，优先改进那些

2.2 Cleaning up incorrectly labeled data

Deep learning are quite robust to random errors in train set.

但对结构性错误就会有问题,

对 dev, test set 上的 incorrectly labeled data, 约 80%

评估就修正，否则就丢弃,

同时修正 dev, test set 保证他们分布相同

2.3 Build your first system quickly, then iterate
尽快建立一个系统，可以知道 bias, variance 在哪，有办法
解决它。

2.4 Training and testing on different distributions.

训练集很大，但与 dev, test set 不太一样

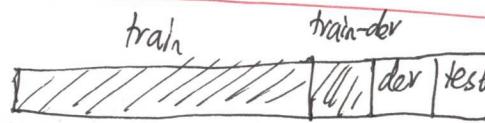
200k 网上图片, 10k 用户上传图片

200k + 5000 为 train set, 2500, 2500 为 dev, test set

2.5 Bias and Variance with mis matched data distributions

若 Train set data 与 dev/test data 未作不同分布，很难确认误差
代表高 Variance 但是高 bias.

把 train set 随机打乱，抽出一个 train-dev set



Train set error: 1% variance

Train-dev error: 9%

Dev error: 10%

Human error --- 0% available bias 0% > bias

Training error: 10%, 10% > bias

Train-dev error: 11%, 11% > mis match

Dev error: 12%, 20% > mis match

2.6 addressing data mismatch

Mismatch 对应歧义，即 differences between training set and dev/test set.

- 一种解决方法是 Artificial data analysis, 但可能只用一个子集去训练 model, 会 overfitting.

2.7 Transfer learning

从一个 model 中取得知识适用于另一个任务中.

什么是 transfer learning: 当 连续 两个 问题 中 有 共同 数据, 可以利用已有数据.

方法:

对很困难的图片训练, 把最后一层拿走, 再放回, 同时
 pre-training 存权重, 对新的数据进行训练, 甚至引入
 fine tuning.

2.8 Multi-task learning

一张图片同时识别 4 个人, 车辆, 停车牌志, 交通灯,

pedestrians	$x^{(i)}$, (4×1)	$y^{(i)}$
Cars	0	$y^{(1)}$
stop sign	1	$y^{(2)}$
traffic light	1	$y^{(3)}$
	0	\dots
		$y^{(m)}$
		$(4 \times m)$

对一个输出 \hat{y} , $2 \cdot 4 \times 1$, 对应于训练集中的平均损失:

$$\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^4 L(\hat{y}_j^{(i)}, y_j^{(i)})$$

$$\text{logistic loss } L(\hat{y}_j^{(i)}, y_j^{(i)}) = -y_j^{(i)} \log \hat{y}_j^{(i)} - (1 - y_j^{(i)}) \log (1 - \hat{y}_j^{(i)})$$

2.9 What is end-to-end deep learning

用单个神经网络训练. (相对单个神经网络好用)

2.10 Whether to use end-to-end learning

优点: ① let the data speaks

② less hard-designing of components needed.

缺点: ① May need large amount of data
 ② Excludes potentially useful hand-designed components.

Course 4, Convolutional Neural Networks

week 1, Foundations of CNN

1.1 Computer Vision

1.2 Edge detection example.

一个 6×6 图像，矩阵为 $6 \times 6 \times 1$.

挑选一个 3×3 的 filter., 算法符号“ $*$ ”

3×3 的框从左上角, 右移, 下移, 直到最右

$$\begin{array}{c} * \quad \begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline -5 & -4 & 0 & 8 \\ \hline -10 & -2 & 2 & 3 \\ \hline 0 & -2 & 4 & -7 \\ \hline -3 & -2 & -3 & -16 \\ \hline \end{array} \end{array}$$

3×3 4×4

用 3×3 filter 做 6×6 图片做 convolution.

输出是一个 4×4 矩阵.

根据①: $\begin{pmatrix} 3 & 0 & 1 \\ 1 & 5 & 8 \\ 2 & 7 & 2 \end{pmatrix} * \begin{pmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{pmatrix} = 3 \times 1 + 1 \times 1 + 2 \times 1 + 0 \times 0 + 5 \times 0 + 7 \times 0 - 1 \times 1 - 8 \times 1 - 2 \times 1 = -5$

先 element-wise products, 再求和。
依次计算。

最左边是 6×6 图片, 中间是 filter, 右边是另一张 4×4 图片

如何上面对应可以用于 vertical edge detection?

$$\begin{array}{ccccccccc} 10 & 10 & 10 & 0 & 0 & 0 & & & \\ 10 & 10 & 10 & 0 & 0 & 0 & & & \\ 10 & 10 & 10 & 0 & 0 & 0 & * & \begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array} & = & 0 & 30 & 30 & 0 \\ 10 & 10 & 10 & 0 & 0 & 0 & & & \\ 10 & 10 & 10 & 0 & 0 & 0 & & & \\ 10 & 10 & 10 & 0 & 0 & 0 & & & \\ 10 & 10 & 10 & 0 & 0 & 0 & & & \\ \end{array}$$

6×6 3×3 4×4

在高亮, 中间元素为垂直边缘
 \rightarrow 垂直边缘

1.3 More Edge detection

上面的输出矩阵左右翻转一下,

$$\begin{array}{ccccccccc} 0 & 0 & 0 & 10 & 10 & 10 & & & \\ 0 & 2 & 0 & 10 & 10 & 10 & * & \begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array} & = & 0 & -30 & -30 & 0 \\ 0 & -2 & 0 & 10 & 10 & 10 & & & \\ 0 & 0 & 0 & 10 & 10 & 10 & & & \\ 0 & 0 & 0 & 10 & 10 & 10 & & & \\ 0 & -30 & -30 & 0 & 0 & 0 & & & \\ 0 & 30 & 30 & 0 & 0 & 0 & & & \\ 0 & 30 & 30 & 0 & 0 & 0 & & & \\ 0 & 30 & 30 & 0 & 0 & 0 & & & \\ 0 & 30 & 30 & 0 & 0 & 0 & & & \\ \end{array}$$

6×6



Vertical Edge detection

$$\begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array}$$

Horizontal Edge detection

$$\begin{array}{|c|c|c|} \hline 1 & 1 & 1 \\ \hline 0 & 0 & 0 \\ \hline -1 & -1 & -1 \\ \hline \end{array}$$

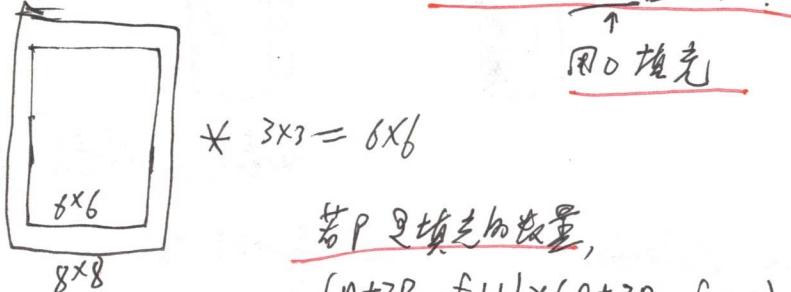
把这九个滤波器合成多长, 通过 back prop, 让 CNN 自己衍生出

1.4 Padding 填充

对 $n \times n$ 图像，用 $f \times f$ 的 filter，输出为 $(n-f+1) \times (n-f+1)$

两个独立的① 固体填充

② 角落边缘的像素又被记录一次，边缘像素
为重复存在。在卷积之前在图像外填充图像。



$$* 3 \times 3 = 6 \times 6$$

若 P 是填充的数量，

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

2种填充方法：一种是 valid 填充，不填充

一种是 same 填充，输出与输入大小一样
填充后，

$$n+2p-f+1 = n \Rightarrow p = \frac{f-1}{2}$$

(stride 步幅)

f 通常为 odd.

1.5 Strided convolutions.

上一节中 stride = 1，即一步，向左却没移一格
还可令 stride = 2, 3, ...

$f \times f$ 的 filter 对 $n \times n$ 的图像，padding 为 p ，

stride 为 s ，输出为 $\left(\frac{n+2p-f}{s}+1\right) \times \left(\frac{n+2p-f}{s}+1\right)$

若 $\frac{n+2p-f}{s}+1$ 不是整数，则向下取整 floor.

但只有填在图像 (or 填之后) 内部计算。

真正的差是先将 filter 翻转， $\begin{bmatrix} 3 & 4 & 5 \\ 1 & 0 & 2 \\ -1 & 9 & 7 \end{bmatrix}$ 变为 $\begin{bmatrix} 7 & 2 & 5 \\ 9 & 0 & 4 \\ -1 & 1 & 3 \end{bmatrix}$
相当于 3×3 的滤波器。

用倒转后的 ~~filter~~ 与 $n \times n$ 图像 element-wise product.

但在 CNN 中，用不翻转的 filter，(且请看 cross-correlation convolution)

1.3 · Convolutions over volumes.

$6 \times 6 \times 3$ 的红色图， $3 \times 3 \times 3$ 的 filter，输出 $4 \times 4 \times 1$

逐层匹配

$$\text{红色通道: } \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix}, \text{ 绿色 } \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix}, \text{ 蓝色 } \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix}$$

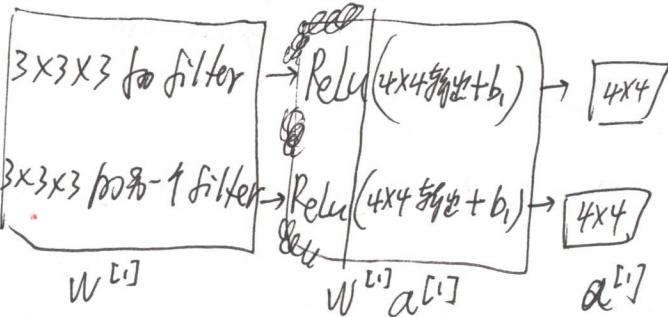
也可用 2 个或者 4 个 $3 \times 3 \times 3$ 的 filter.

维度： $n \times n \times n_c$ (通道数) 图像， $f \times f \times n_c$ filter，

输出： $(n-f+1) \times (n-f+1) \times n_c$ ， n_c 为用的 filter 个数

1.7 One layer of a convolutional network.

$6 \times 6 \times 3$ input
 $\alpha^{[0]}$



$$z^{[1]} = W^{[1]} \alpha^{[0]} + b^{[1]}$$

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

$f^{[l]}$ = filter size

$p^{[l]}$ = padding

$s^{[l]}$ = stride

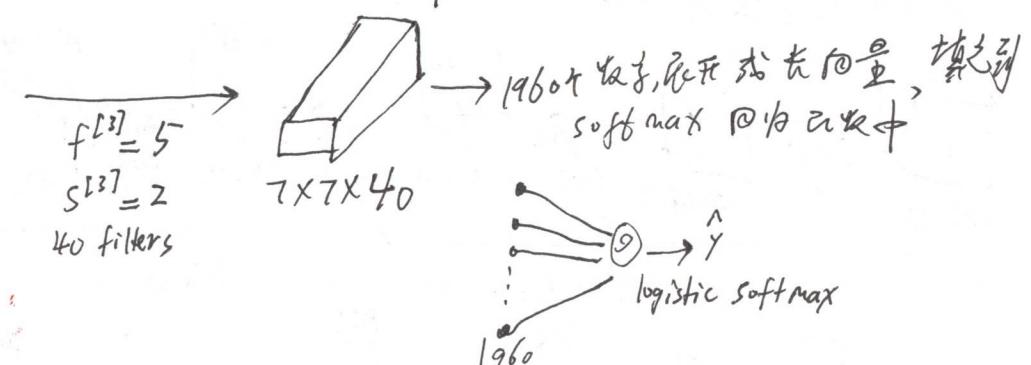
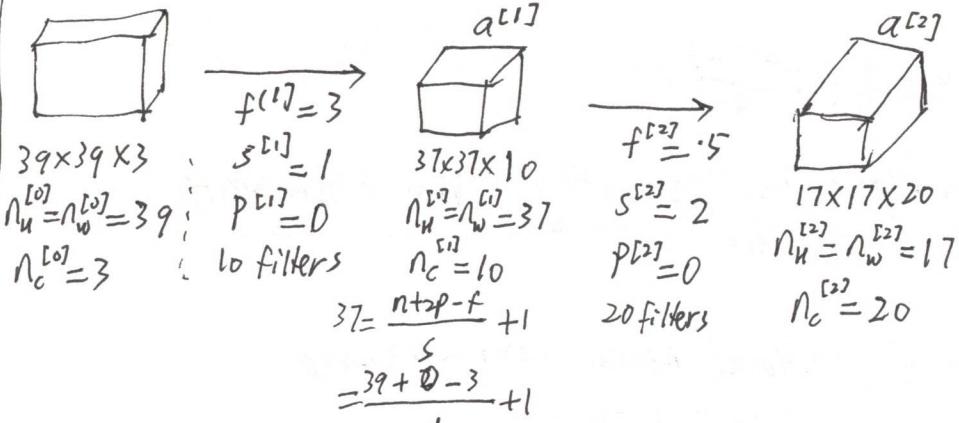
Each filter: $f^{[l]} \times f^{[l]} \times n_c^{[l+1]}$

Activation: $\alpha^{[l]} \rightarrow n_h^{[l]} \times n_w^{[l]} \times n_c^{[l]}$

Weights: $A^{[l]} = m \times n_h^{[l]} \times n_w^{[l]} \times n_c^{[l]}$

bias: $n_c^{[l]} \rightarrow (1, 1, 1, n_c^{[l]})$

1.8 A simple convolution network example.

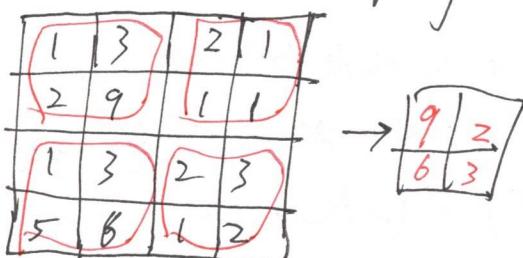


Types of layer in a convolutional network:

- (1) Convolution (CONV)
- (2) Pooling (POOL)
- (3) Fully connected (FC)

1.9 Pooling layer

{3/2: pool layer: max pooling



执行 max pool to 2x2
 将其与之并列，输出层的尺寸是
 1/4.

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

但有特殊情况 $0 - 0 \leq 0$ 成立,
为了避免这种情况,

$$\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha \leq 0$$

\uparrow
margin

loss func:

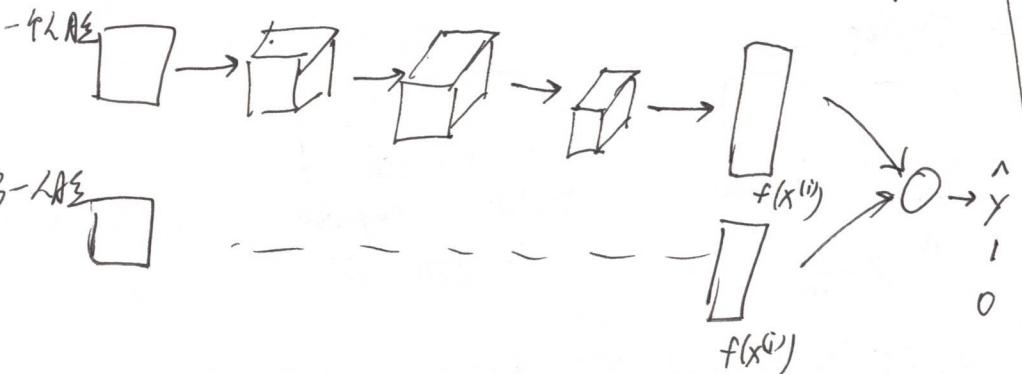
$$L(A, P, N) = \max(\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha, 0)$$

10000 pictures of 1000 persons.

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

选择对称训练的三元组, where $d(A, P) \approx d(A, N)$

4.5 Face verification and binary classification.



Siamese 网络, 将待识别的 128 维特征输入到网络①为单向的元, 因此
没有, 相同为 1, 否则为 0

$$\hat{y} = \sigma \left(\sum_{k=1}^{128} w_i |f(x^{(i)})_k - f(x^{(j)})_k| + b \right)$$

\downarrow 代表这个向量中的第 k 个元素.

$$\text{其他形式: } \chi^2 \text{ 公式: } \frac{(f(x^{(i)})_k - f(x^{(j)})_k)^2}{f(x^{(i)})_k + f(x^{(j)})_k}$$

首先算好每个类别的编码, 当识别时, 用卷积网将计算
编码与现存的编码比较, 进行预测.

4.6 What is Neural net style transfer?

4.7 What are deep convNet learning

深层生成网训练, $\square \rightarrow \square$, 从 texture, 到更复杂的物体

4.8 Cost func.

包含 C , 风格 S , 特征的网 G .

评估生成图片的损失: 通过 cost func.

$J_{content} \leftarrow J_{content}(C, G), J_{style}(S, G)$

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

2.3 · Residual Network (ResNets)

skip connection: “层反馈信号另外一层 (skip layer)”

residual block:



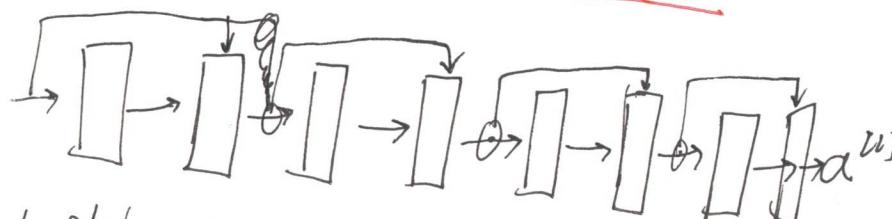
✓ $a^{[l]}$ → Linear → ReLU $a^{[l+1]}$ → Linear → ReLU → $a^{[l+2]}$

$a^{[l]}$ → Linear → ReLU $a^{[l+1]}$ → skip connection → $a^{[l+2]}$

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

$\Delta a^{[l+2]} = g(z^{[l+2]} + a^{[l]})$

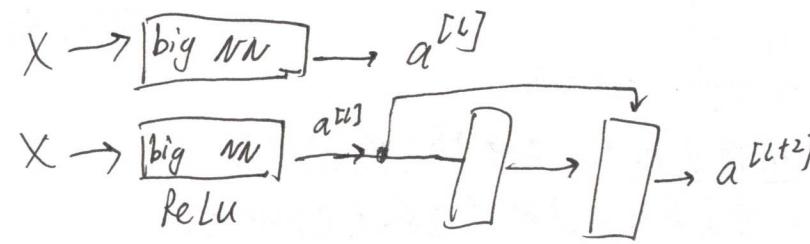
in ResNet, $a^{[l]}$ 直接向后, 携带到网络的下一层
④ Residual block 可以将更平的 NN.



在 plain network 上加上 skip connection 就是 ResNet.

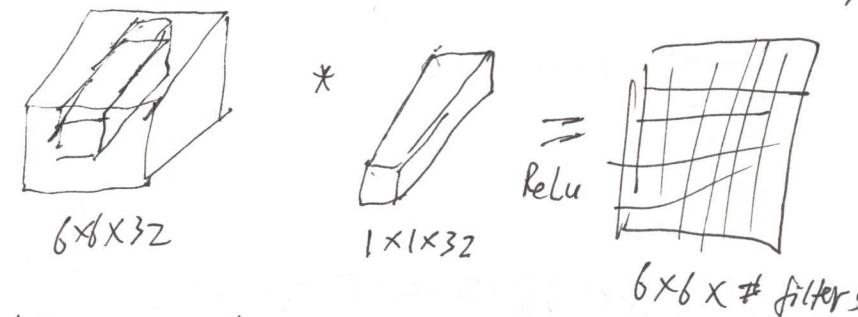
2.4 why ResNets work?

再去好好理解 why Res Network?
回头看 code.



$$\begin{aligned} a^{[l+2]} &= g(z^{[l+2]} + a^{[l]}) \\ &= g(W^{[l+1]} a^{[l]} + b^{[l+2]} + a^{[l]}) \\ &\text{if } W^{[l+2]} = 0, b^{[l+2]} = 0, a^{[l+2]} = g(a^{[l]}) = a^{[l]} \end{aligned}$$

2.5 Network in Network and 1×1 convolution



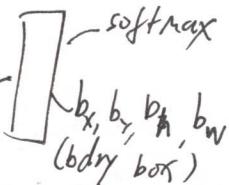
$1 \times 1 \times 32$ 的 32 表样输出: 一个特征元输入是 32 个数字, 而相隔在
和宽度上每一个片上的 32 个数字, 32 个数字有不同的通道,
即 32 个权重, 用 ReLU 为活性 func, 选出该重.
该方法能将 1×1 convolution or network in network,

它将 NN 增加一个 nonlinear func, 从而 ↓ 或保持输入层中的通道
因子不变. TBM

Week 3 Object detection:

3.1 Object localization.

图 R → convolution



Refining target label

1. pedestrian
2. car
3. motorcycle
4. background

$$y = \begin{cases} p_c \\ b_x \\ b_y \\ b_t \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{cases}$$

总共有 1~3 及 4, 共 5 个类



loss func:

$$L(\hat{y}, y) = (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots + (\hat{y}_8 - y_8)^2$$

3.2 Landmark detection.

Q. 3.3 Object detection

用滑动窗口的方法. (sliding window detection)

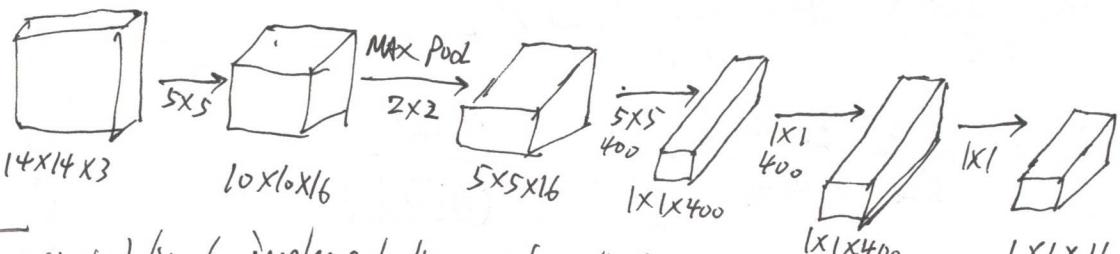
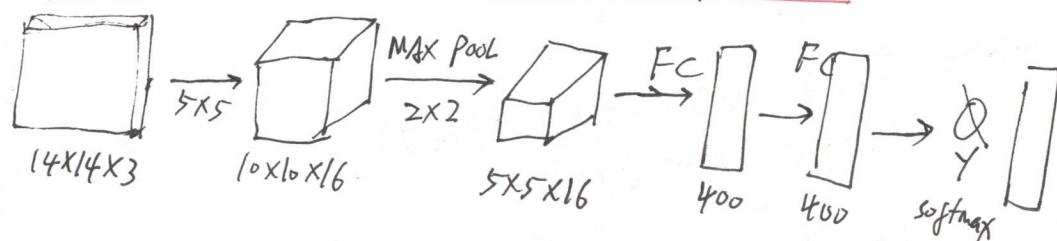
一个子窗口, 从左上角, →, ↓ 移动

大一点的窗口, -----.

再大一点. -----.

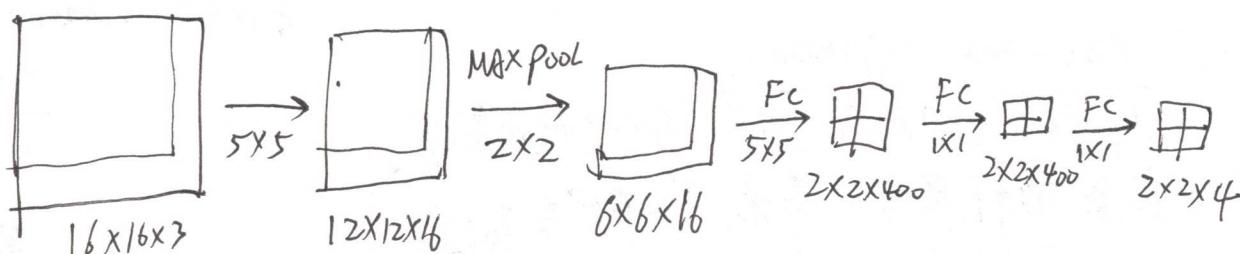
3.4 Convolutional implementation of sliding windows.

Turning FC layer into convolutional layers



Convolutional implementation of sliding windows

输入图 R 为 16x16x3, 经 CNN 后为 14x14x3 图片



3.8 Attention Model

插上反图：

$$a^{<t>} = (\vec{a}^{<t>}, \vec{\alpha}^{<t>})$$

$$\sum_{t'} \alpha^{<1, t'} = 1$$

$$C^{<1>} = \sum_{t'} \alpha^{<1, t'} a^{<t>}$$

$\alpha^{<t, t'>} = \text{amount of attention should } y^{<t>} \text{ pay to } a^{<t'>}$

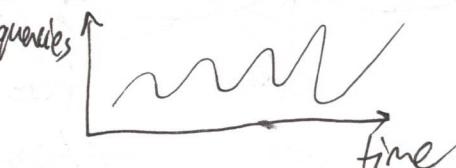
$$C^{<2>} = \sum_{t'} \alpha^{<2, t'} a^{<t'>}$$

$$\alpha^{<t, t'>} = \frac{\exp(e^{<t, t'>})}{\sum_{t'=1}^T \exp(e^{<t, t'>})}$$

$$S^{<t-1>} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow e^{<t, t'>}$$

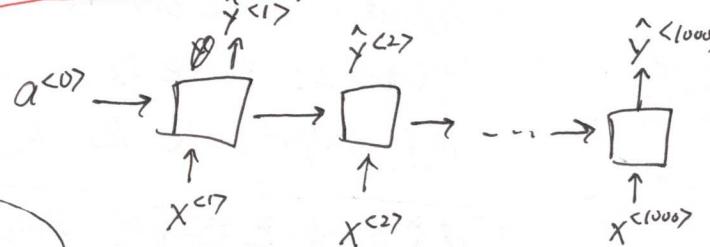
q - speech recognition

audio $\in \mathbb{R}^{T \times F}$ to spectrogram.



CTC (connectionist temporal classification)

CTC cost for speech recognition.



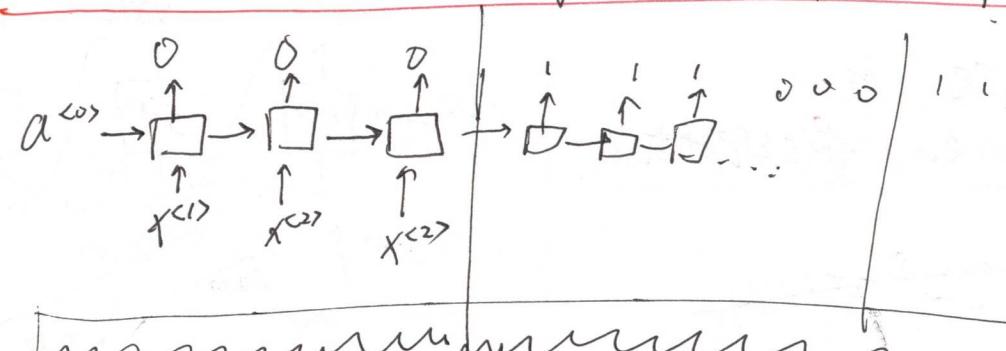
"the brown quick brown fox"

t t t - . h - e e e - - - l u - - - g g g - - space. \rightarrow the q

CTC cost function: 把空白符之间的重叠的字符拼接起来。

3.10 Trigger word detection.

在音频流中 trigger word 之差为 0, 触发词检测 |



3.9 ~~YOLO~~ put it together: YOLO algorithm.

3.10 Region proposals

R-CNN: 带 RPN 的 CNN.

不是一维滑动窗口，而是选择某些窗口，
这些窗口基于图像分割算法。

Faster R-CNN

Faster R-CNN

Week 4 Special applications: Face recognition & Neural style transfer.

4.1 What is face recognition.

Face verification vs. face recognition.

Verification:

- input image, name, ID

- output whether the input image is that of the claimed person

Recognition:

- has a database of K persons

- get an output image

- output ID if the image is any of the K persons

4.2 One-shot learning

Learn from one example to recognize the person again.

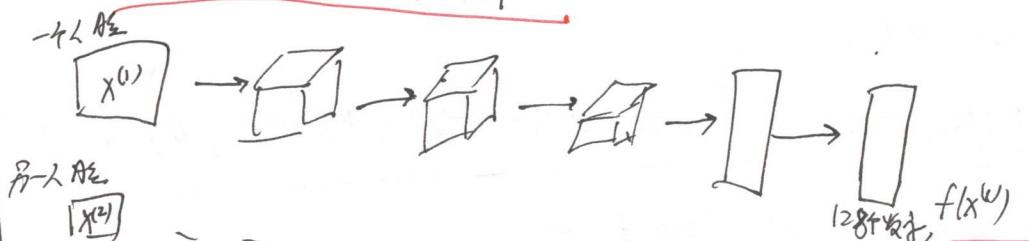
Learning a "similarity" func.

$d(\text{img}1, \text{img}2) = \text{degree of difference between images}$

If $d(\text{img}1, \text{img}2) \leq t$ "same"

$> t$ "different"

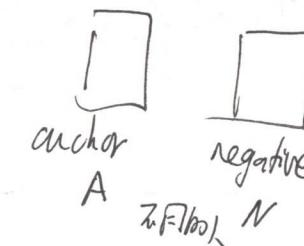
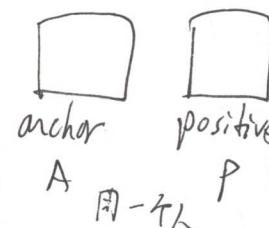
4.3 Siamese network



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

in the siamese network $x^{(1)}$, $x^{(2)}$ is the same person, then $d(x^{(1)}, x^{(2)})$ is small.

4.4 Triplet Loss



triplet means A, P, N 三张图片。

Find the generated image G

1. Initiate G randomly, $G: 100 \times 100 \times 3$
2. Use gradient descent to minimize $J(G)$

$$G := G - \frac{\partial}{\partial G} J(G)$$

4.9 Content Cost func:

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

$J_{content}(C, G)$:

use hidden layer l to compute $J_{content}(C, G)$

先让 G 不同于 C , 用一个预训练好的卷积模型

- let $a^{[l](C)}$ and $a^{[l](G)}$ (VGG or other).
 l on the images be the activation of layer

- If $a^{[l](C)}$ and $a^{[l](G)}$ are similar, both images have similar content.

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

4.10 Style cost func.

use layer l 's activation to measure 'style'.

Use activation:

卷积层不同通道的

activation 为 correlate ？

13.6 to 54 channels



let $a_{i,j,k}^{[l]}$ = activation at (i, j, k)

$G^{[l](s)}$ is $n_c^{[l]} \times n_c^{[l]}$, n_c 为通道数

$$G_{kk'}^{[l](s)} = \sum_{i=1}^{n_h^{[l]}} \sum_{j=1}^{n_w^{[l]}} a_{ijk}^{[l](s)} a_{ijk'}^{[l](s)}, \quad k, k' \text{ 为通道 } k, k'$$

$$G_{kk'}^{[l](a)} = \sum_{i=1}^{n_h^{[l]}} \sum_{j=1}^{n_w^{[l]}} a_{ijk}^{[l](a)} a_{ijk'}^{[l](a)} \quad k = 1 \dots n_c$$

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

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

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

4.11 1D and 3D generalizations of models

Course 5, Sequence Models

Week 1 Recurrent NN.

1.1 why Sequence Models?

1.2 Notation

~~从哪到哪的字符的集合？~~

"Harry Potter and Hermione Granger invented a new spell."
 $x^{<1>} x^{<2>} x^{<3>} \dots$

$x^{<t>} \rightarrow$ 在序列中第 t 位置 (输入)
 $y^{<1>} - - - - - (输出)$

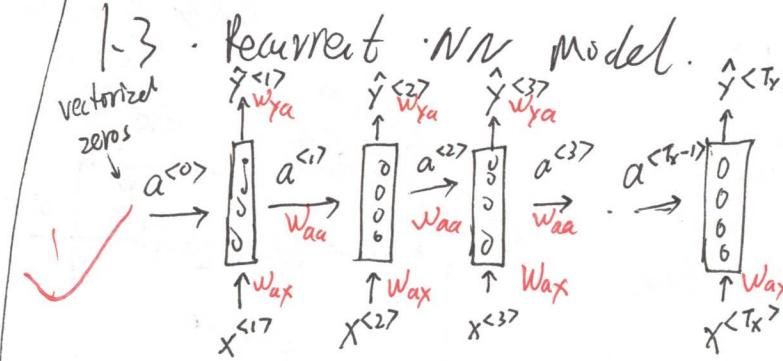
T_x : 输入序列
 T_y : 输出 —

$x^{(i)<t>} = i$ 个样本；中第 i 个元素， $T_x^{(i)}$: 第 i 个输入
 $y^{(i)<t>} = - - - - - \rightarrow T_y^{(i)}$: 第 i 个输出

对于一个 1M 甚至几万、上百万的常用字典，
用 one-hot 表示每个单词。

$x^{<1>} = "Harry"$, 在字典中 4075 位, 故 $x^{<1>} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \leftarrow 4075/1$
上百万个单词有 97 one-hot 向量

若字典里没有的词，通过 unknown word 的方法， $\langle \text{unk} \rangle \notin$



$x^{<1>} \rightarrow$ 从第 1 个 hidden layer, 得到 $y^{<1>} \rightarrow$

$x^{<2>} \rightarrow$ 从第 2 个 hidden layer, 得到 $y^{<2>} \rightarrow$ 依此类推 $a^{<1>} \rightarrow$

一个特点：只使用了直到前面的信息.

$$a^{<0>} = \vec{0}, \quad [a^{<1>} = g_1(w_{aa}a^{<0>} + w_{ax}x^{<1>} + b_a) \leftarrow \text{tanh, ReLU} \\ \hat{y}^{<1>} = g_2(w_{ya}a^{<1>} + b_y) \leftarrow \text{sigmoid, --}$$

$$\begin{cases} a^{<t>} = g(w_{aa}a^{<t-1>} + w_{ax}x^{<t>} + b_a) \\ \hat{y}^{<t>} = g(w_{ya}a^{<t>} + b_y) \end{cases}$$

simplified RNN notation:

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

$$100 \updownarrow \begin{bmatrix} W_{aa} & W_{ax} \end{bmatrix} = W_a \quad (100, 100|100) \\ \updownarrow \updownarrow \quad \updownarrow \updownarrow$$

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

$$[W_{aa}, W_{ax}, W_{ya}, b_a] \updownarrow \quad \updownarrow \quad \updownarrow \quad \updownarrow$$

1.5 Different types of RNN

并不是所有的RNN都是 $T_x = T_y$
many to many

一个常见的例子：

输入一个句子，输出一个句子， many to one.

还有 one-to-one, one-to-many, many-to-many ($T_x \neq T_y$).

1.6 Language model and sequence generation

两个很重要的应用，language model 和 sequence generation
是学习一种语言的频率，

Language modeling with an RNN:

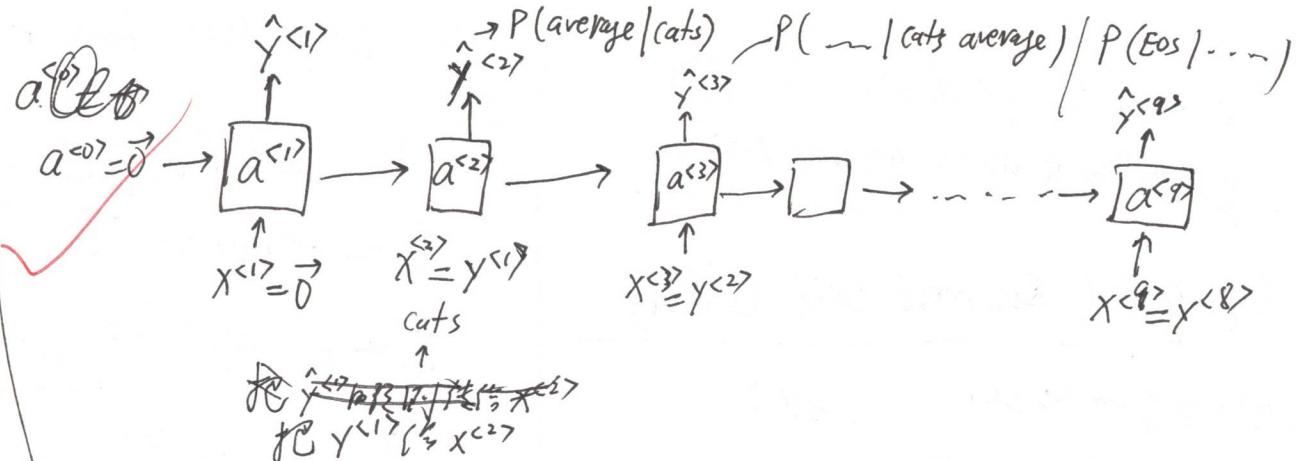
例句：Cats average 15 hours of sleep a day. <EOS>

The Egyptian ~~Mau~~ is a bread of cat. <EOS>
<UNK>

对于<BOS>识别句子模型

没有的词用<UNK>代替

~~输出~~ 打造完一个one-hot，开始进RNN。
输出也是通过神经模型

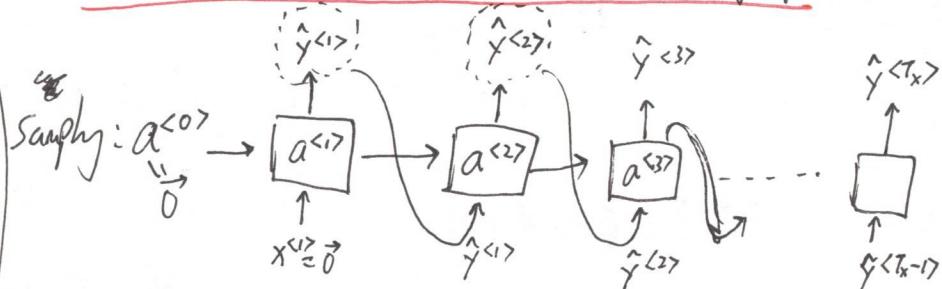


$$\text{Loss func: } L(\hat{y}^{<t>}, y^{<t>}) = -\sum_i y_i^{<t>} \log \hat{y}_i^{<t>}$$

$$\text{总 loss func: } L = \sum_t L^{<t>}(\hat{y}^{<t>}, y^{<t>})$$

1.7 Sample novel sequences.

语言模型学到了什么，做 sampling



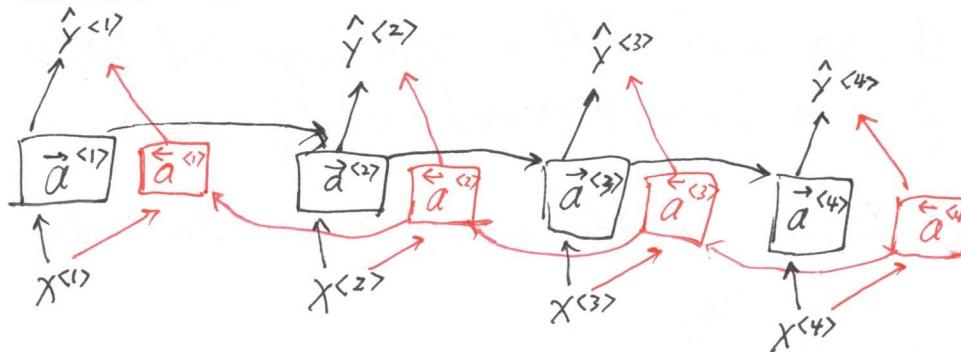
第一步 对语言模型生成的第一句话，输入 $x^{<0>} = 0$, $a^{<0>} = 0$
得到 softmax 指向的频率，根据频率进行随机选择。

把上一步选择得到的 $\hat{y}^{<1>}$ 输入。

之后的每一个步骤都如此。

1.11 Bi-directional RNN

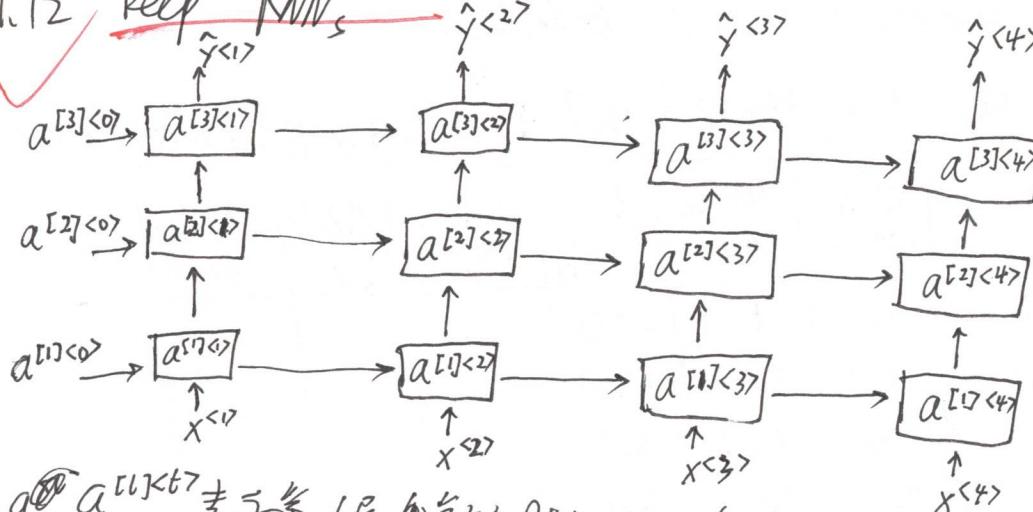
不仅获取前向的信息，还有本素的信息。



$$\text{预测 } \hat{y}^{<t>} = g(W_g[\vec{a}^{<t>}, \vec{a}^{<t>}] + b_y)$$

基本单元可以是 RNN, GRU, LSTM 等

1.12 Recurrent MN



$a^{[t]} < t >$ 表示第 t 层的隐藏状态，由 σ 激活函数激活。

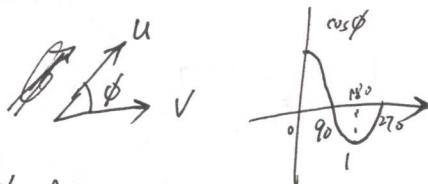
$$a^{[2]} < 3 > = g(W_a^{[2]} [a^{[2]} < 2 >, a^{[1]} < 3 >] + b_a^{[2]})$$

当男性被问及 man 和 woman 所属于 King 和 Queen，
女性被问及 man - woman，找出。

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

找 word w: $\arg \max \text{Sim}(e_w, e_{\text{king}} - e_{\text{man}} + e_{\text{woman}})$

余弦相似度: $\text{sim}(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2} = \cos \phi$

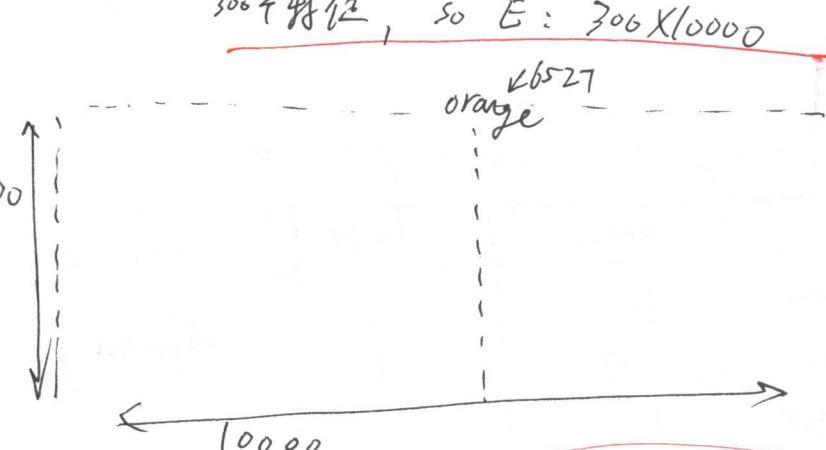


当 $\phi=0^\circ$ 时, $\cos \phi=1$, u, v 共线。

2.4 Embedding matrix

假设一个 embedding matrix E, 10000 个单词，每词

300 维特征，那么 $E: 300 \times 10000$

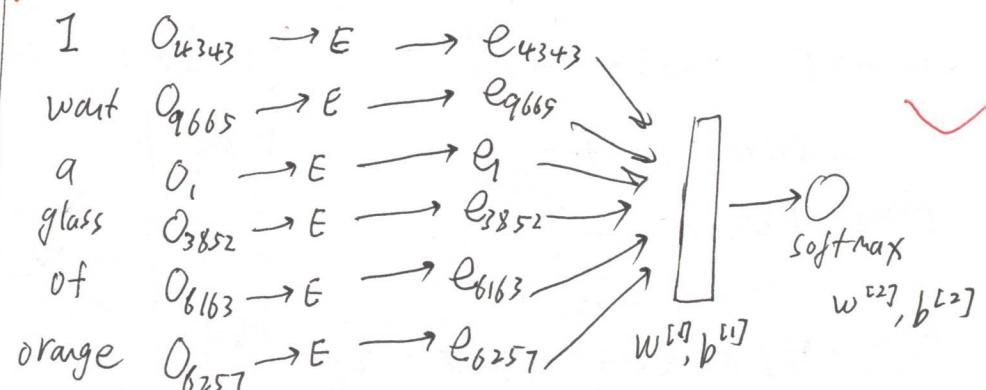


$$E * \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}_{10000 \times 1} = e_{6527}$$

任意单词 w, e_w 在其嵌入向量， O_j 只有第 j 位是 1 为 one-hot 向量。

2.5 Learning word Embedding

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



更常见的是固定权重：

即各单词的向量都是从 4 个单词向量下-一个单词

由之语言模型：用附近的前几个单词预测
上一个

预测嵌入：其他类型上下文

- 4 words on left & right
- Last 1 word
- nearby 1 word

small dataset, $k = 5 \sim 20$
large data set: $k = 2 \sim 5$

Model:

$$\text{Softmax: } P(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10000} e^{\theta_j^T e_c}}$$

$$P(Y=1|C, t) = \sigma(\theta_t^T e_c)$$

Orange
6257

$$\theta_{6257} \rightarrow E \rightarrow e_{6257} \rightarrow \begin{matrix} \rightarrow \text{o juice} \\ \downarrow \rightarrow \text{o king} \end{matrix}$$

如何选负样本?

$$P(w_i) = \frac{f(w_i)^{\frac{3}{4}}}{\sum_{j=1}^{10000} f(w_j)^{\frac{3}{4}}}$$

$f(w_i)$ 为词向量到词在词库中第 i 个英文词的词频

2.8 GloVe Word vectors

(global vectors for word representation)

	X	Y	
	context	word	target
orange	juice	1	
orange	key	0	
orange	bark	0	
orange	the	0	
orange	of	0	
c	↑	↑	
t	↑	↑	
y	↑	↑	

1个正样本 orange
k个负样本

X_{ij} : 单词 i 在单词 j context 中出现的词频

对于 glove, Glove, 之 $\xrightarrow{\text{context}}$ context 和 target 之间的互换位置相近的词.

Glove 和 glove 之间差距是 1/6:

minimize $\sum_{i=1}^{10000} \sum_{j=1}^{10000} f(X_{ij})(\theta_i^T e_j + b_i + b_j - \log X_{ij})^2$

当 $X_{ij} = 0$, $\log X_{ij}$ 为无穷大, 故有 \rightarrow 加权 $f(X_{ij})$, $f(0) = 0$

$f(X_{ij})$ 用于对 of, the 这些高频无意义词的权重, 以及罕见不常见词的权重.

$$(A\theta_i)^T (A^T \theta_j) = \theta_i^T A^T A^{-T} \cdot \theta_j = \theta_i^T \theta_j$$

2.9 sentiment classification

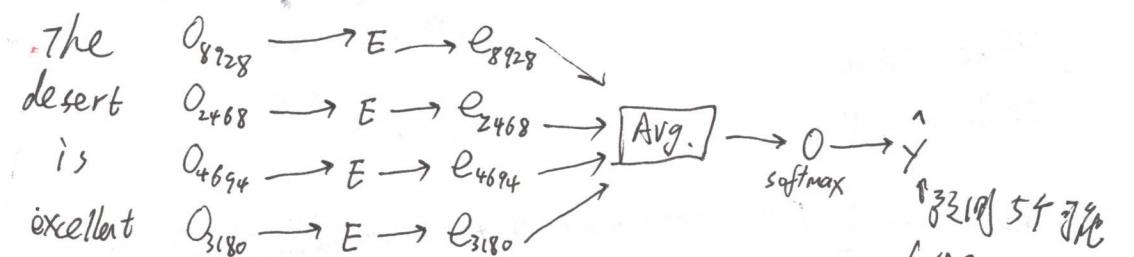
"The desert is excellent"

8928

2468

4694

3180



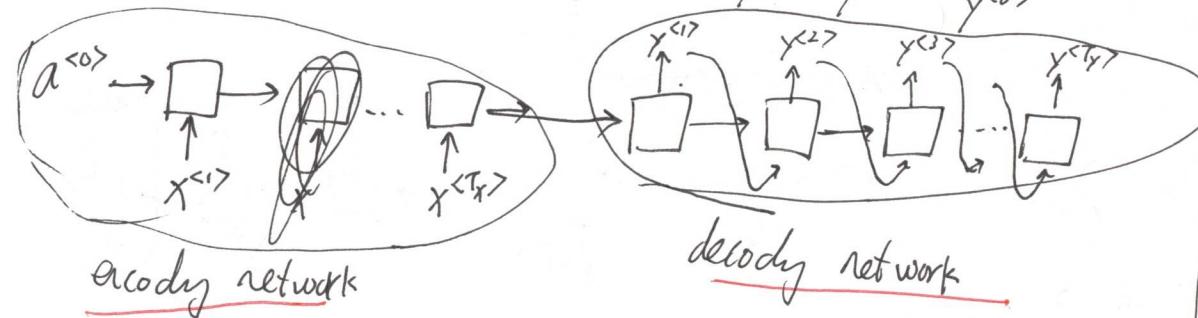
缺点: 特异名词.

Week 3 · Sequence Models & Attention mechanism

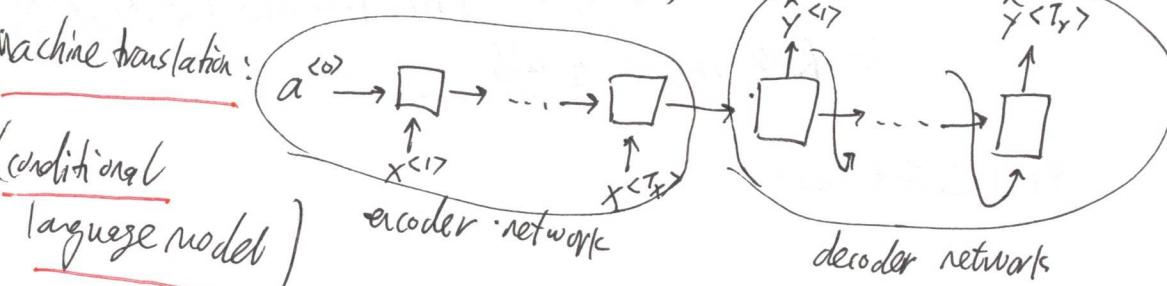
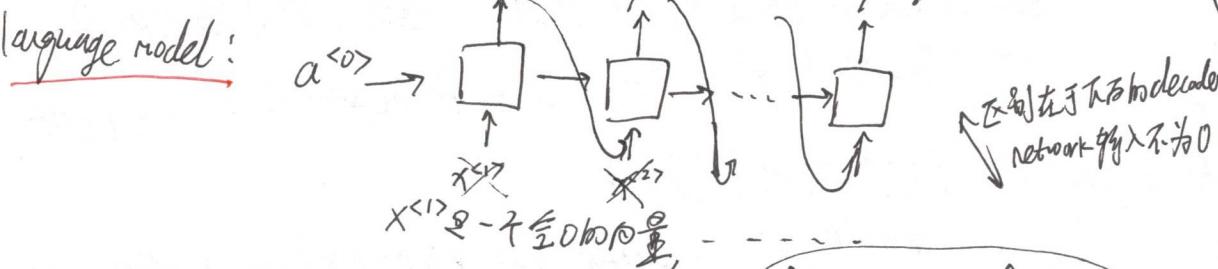
1.1 Various sequence to sequence architectures

输入法语:
输出英语:

Jane visite l'Afrique en septembre
 Jane is visiting Africa in September.



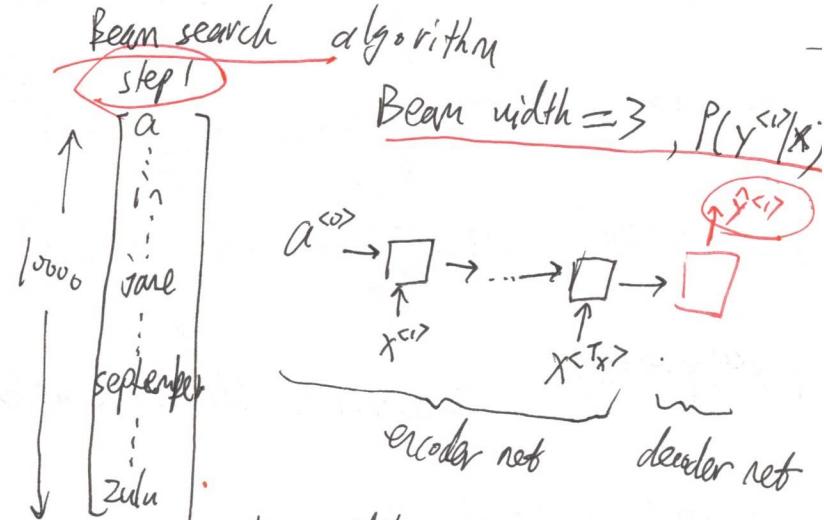
3.2 Picking the most likely sentence



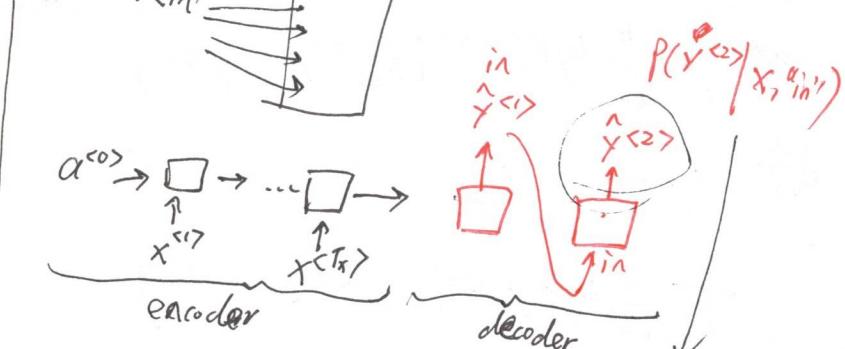
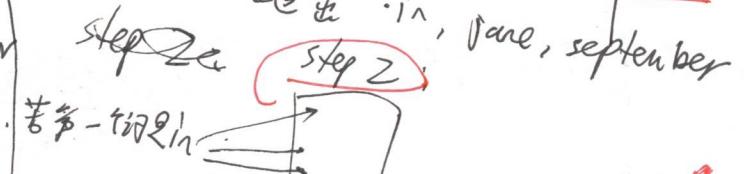
$$\arg \max_{y^{<1>} \dots y^{<T_y>}} p(y^{<1>} \dots y^{<T_y>} | x)$$

用 beam search 找出合适的 y 使 $p(y|x)$ 最大。
 Greedy search 从第一个词搜索，这里更需要整句的搜索。

3.3 Beam search



beam width = 3 means 第一个阶段有 3 个
 可能的候选句：Jane, in, september.



$$P(y^{<1>} \dots y^{<T_y>} | x) = P(y^{<1>} | x) P(y^{<2>} | x, y^{<1>}) \dots P(y^{<T_y>} | x, y^{<1>} \dots y^{<T_y-1>})$$

3.8 Bleu Score

Bilingual Evaluation under study.

对于两个不同的翻译，BLEU给出评分；

一个具体的例子：

翻译输出：the the the the the the the

参考翻译1：The cat is on the mat.

precision = $\frac{7}{7}$

虽然不好，改进：

修改1，the出现2次，修改2，the出现2次
故上限为2。

precision = $\frac{2}{7}$

且 P_n 为单词the的总出现次数

分子是单词the的出现次数，但到达上限时就断。

BLEU score on bigrams

$$P_n = \frac{\sum_{n\text{-grams}}^{} \text{count}_{clip}(n\text{-gram})}{\sum_{n\text{-grams}}^{} \text{count}(n\text{-gram})}$$

$$\text{BLEU}(\%) = \exp\left(\frac{1}{4} \sum_{n=1}^4 P_n\right)$$

Evaluating machine translation:

French: Le chat est sur le tapis

Reference1: The cat is on the mat.

Reference2: There is a cat on the mat.

MT. output: the the the the the the

precision: $\frac{7}{7}$

modified precision: $\frac{2}{7}$

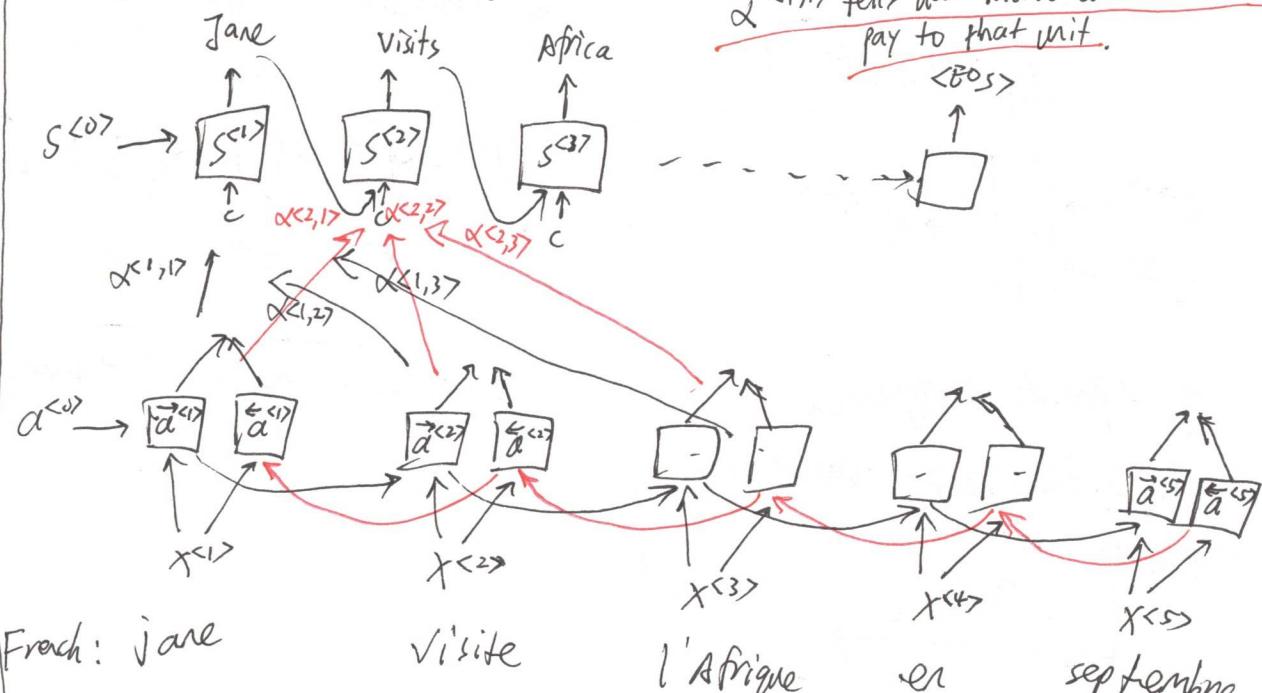
通常会加上惩罚项 Brevity penalty, BP, $BP \cdot \exp\left(\frac{1}{4} \sum_{n=1}^4 P_n\right)$

$$BP = \begin{cases} 1 & , \text{ if } \text{MT_output_length} > \text{reference_output_length} \\ \exp(1 - \text{MT_output_length}/\text{reference_output_length}) & , \text{ otherwise} \end{cases}$$

3.7 Attention Model Intuition

对于这个 encoder-decoder 模型，Bleu score 100%.

Attention model intuition:



French: Jane

visits

l'Afrique

er

septembre