

1. Training set size used in gradient descent: Batch, mini-batch, tochastic gradient descenet

$W := W - \alpha \cdot dw$
 $dw = 1/m \cdot \sum dw$ (under $x(i)$):

Cost function J is the sum of every training example error-->
 Gradient descent is to make this on every example small-->minimize the sum of every training example error

$dw(i)$: parameter gradient descent direction, for making only $J(i)$ -example $x(i)$ cost smaller, but could not make sure sum of $J(i)$ -all examples' error, is smaller.
 ----> $dw(i)$ average over m size training example, make sure $J(i)$ over m , is becoming smaller.

Logistic regression on m examples

云课堂

$$J = 0; \quad dw_1 = 0; \quad dw_2 = 0; \quad db = 0$$

For $i = 1$ to m

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

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

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

$$\frac{\partial J}{\partial z^{(i)}} = a^{(i)} - y^{(i)}$$

$$\frac{\partial J}{\partial w_1} += x_1^{(i)} \frac{\partial J}{\partial z^{(i)}} \quad \frac{\partial J}{\partial w_2} += x_2^{(i)} \frac{\partial J}{\partial z^{(i)}} \quad \frac{\partial J}{\partial b} += \frac{\partial J}{\partial z^{(i)}} \quad \downarrow n=2$$

$J = m$
 $w_1 := m; \quad dw_1 := m; \quad db := m.$

Variation



20人赞同了该回答

因为优化是一门独立的学科，梯度下降算法本身就不是为机器学习的问题而设计（梯度下降最早可追溯到200多年前的柯西），在梯度下降的大部分应用中，并没有“样本”的概念存在。

梯度下降的应用问题一般就是抽象为就是简单的最小化一个函数可导函数： $\min_x f(x)$ ，对于这个问题，梯度下降就是很简单的迭代： $x_{t+1} = x_t - \eta \nabla f(x_t)$ 。

只不过在机器学习中，我们有了很多个样本，而目标函数往往就是在每个样本上的loss求平均，也就是说此时问题有了额外的结构： $f(x) = \frac{1}{n} \sum_{i=1}^n \ell_i(x)$ 。此时梯度下降算法就自然变成了：

$$x_{t+1} = x_t - \eta \sum_{i=1}^n \nabla \ell_i(x_t)$$

于是就有了你所说的“梯度求和”的形式。

当然，虽然梯度下降在求解机器学习问题中是一个有效的算法，但因为机器学习问题有更加特殊的结果，那么自然而然的大家就开始思考能否利用这种结构设计出更加高效的算法，于是就有了随机梯度下降算法（SGD）以及它的众多变种，其核心思想与你的想法有些类似，每次只算一个 $\ell_i(x)$ 的梯度，SGD中这个 i 是随机抽样得到的。

至于每轮迭代用最大的那个样本梯度显然是不可行的。一个很容易可以想到的情况就是，如果有一个最大的样本梯度是朝一个方向的，但剩下所有的样本梯度都是朝另一个方向的，你觉得此时应该按最大的那个，还是按剩下所有的？一个更具体的例子，假设 $\ell_1(x) = (x+2)^2$ ，但 $\ell_2(x) = \dots = \ell_{10}(x) = (x-1)^2$ ，很容易算出整个 $f(x)$ 的最小值在 $x = 0.7$ 的位置取到。那么假设我现在在 $x_0 = 0$ 出发想迭代求解这个问题，我应不应该按 $-\nabla \ell_1(x_0)$ 所指示的往负数方向走呢？

发布于 2017-11-18

2. Logistic cost function:

Linear cost function: distance: $(h-y)^2$;
 Logistic regression cost function:
 distance of predicted value and label value is good way to measure prediction performance/error.
 Yet only for using gradient descent to minimize training error, need to make sure cost function J is convex.
 Logistic regression distance error is not convex, that's why choose log type-convex.
 (squared error for sigmoid function-->could not guarantee updating parameter when error/cost is big
 -->converge slow:
 $dw = (y^{\wedge} - y) \cdot (y^{\wedge} (1 - y^{\wedge})) \cdot X^T$
 -->cross cost function: $dw = (y^{\wedge} - y) \cdot X^T$

Logistic Regression cost function

云课堂

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

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

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

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

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

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

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

3. training data -pre-processing

for color image, have 3 channel, -->feature vector: unrolled image into $(w \times h \times 3, 1)$ vector.

Example: Cat vs Non-Cat

The goal is to train a classifier that the input is an image represented by a feature vector, x , and predicts whether the corresponding label y is 1 or 0. In this case, whether this is a cat image (1) or a non-cat image (0).



An image is store in the computer in three separate matrices corresponding to the Red, Green, and Blue color channels of the image. The three matrices have the same size as the image, for example, the resolution of the cat image is 64 pixels X 64 pixels, the three matrices (RGB) are 64 X 64 each.

The value in a cell represents the pixel intensity which will be used to create a feature vector of n -dimension. In pattern recognition and machine learning, a feature vector represents an object, in this case, a cat or no cat.

Forward: $Z = W^T \cdot X + b$; $A = g(Z)$ (1, m)

Backward: $dZ = A - Y$, (1, m)
 $db = 1/m \cdot \sum dp$, Sum (dZ) (1,1)
 dw (n,1) = $1/m \cdot X \cdot dZ^T$.

Note: if need many iterations, for loops still needed:
 for 1: n :
 Forward;
 Backward;
 Gradient descent.

To create a feature vector, x , the pixel intensity values will be "unroll" or "reshape" for each color. The dimension of the input feature vector x is $n_x = 64 \times 64 \times 3 = 12288$.

$$x = \begin{bmatrix} 255 \\ 231 \\ 42 \\ \vdots \\ 255 \\ 134 \\ 202 \\ \vdots \\ 255 \\ 134 \\ 93 \\ \vdots \end{bmatrix}$$

red
green
blue

Implementing Logistic Regression

云课堂

$$J = 0; \quad dw_1 = 0; \quad dw_2 = 0; \quad db = 0$$

For $i = 1$ to m :

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

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

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

$$\frac{\partial J}{\partial z^{(i)}} = a^{(i)} - y^{(i)}$$

$$\frac{\partial J}{\partial w_1} += x_1^{(i)} \frac{\partial J}{\partial z^{(i)}} \quad \frac{\partial J}{\partial w_2} += x_2^{(i)} \frac{\partial J}{\partial z^{(i)}} \quad \frac{\partial J}{\partial b} += \frac{\partial J}{\partial z^{(i)}}$$

$$J = m; \quad dw_1 := m; \quad dw_2 := m; \quad db := m.$$

5. Broadcasting:

making algorithm faster in python
 sum in matrix column: $cal = A \cdot \text{sum}(\text{axis}=0)$; $\text{axis}=1$: vertical computation
 sum in matrix row: $\text{percentatage} = A / \text{cal} \cdot \text{reshape}(1,4)$; each A column divided by corresponding cal column value

adding: $A(m, n) + B(1, n)$: B will be multiplied row to (m, n) matrix, then add with A .
 adding: $A(m, n) + B(m, 1)$: B will multiply columns to (m, n) matrix, then add A .

General Principle

$$\begin{bmatrix} m & n \end{bmatrix} + \begin{bmatrix} 1 & n \end{bmatrix} \rightarrow \begin{bmatrix} m & n \end{bmatrix}$$

$$\begin{bmatrix} m & n \end{bmatrix} + \begin{bmatrix} m & 1 \end{bmatrix} \rightarrow \begin{bmatrix} m & n \end{bmatrix}$$

$$\begin{bmatrix} m & 1 \end{bmatrix} + \begin{bmatrix} 1 & n \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix} + \begin{bmatrix} 100 & 100 & 100 & 100 \end{bmatrix} = \begin{bmatrix} 101 & 101 & 101 & 101 \end{bmatrix}$$

Matrix/Column broadcast

Broadcasting example

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

	Apples	Beef	Eggs	Potatoes
Carb	56.0	0.0	4.4	68.0
Protein	1.2	104.0	52.0	8.0
Fat	1.8	135.0	99.0	0.9

$\text{cal} = A \cdot \text{sum}(\text{axis}=0)$
 $\text{percentatage} = 100 \cdot \text{cal} / \text{cal} \cdot \text{reshape}(1,4)$

Broadcasting example

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 100 & 100 \end{bmatrix} = \begin{bmatrix} 101 & 102 & 103 \\ 104 & 105 & 106 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 100 & 100 \\ 200 & 200 & 200 \end{bmatrix} = \begin{bmatrix} 101 & 102 & 103 \\ 304 & 305 & 306 \end{bmatrix}$$

6. Tech: matrix in Python

>1. List column and row number, instead of only one -in this case python will take this value as row and column both in later computation.
Esp. when creating vector: do not use rank 1 array.

>could use 'assert' to make sure vector dimension;
>Cold use 'reshape' to adjust matrix dimension.

Python/numpy vectors

```
a = np.random.randn(5)
a.shape = (5,)
"rank 1 array"
```

Don't use

```
a = np.random.randn(5,1) → a.shape = (5,1) column vector ✓
```

```
a = np.random.randn(1,5) → a.shape = (1,5) row vector ✓
```

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

7. Logistic regression cost function

>1. **Principle**: based on the label training data (x, y) , find model parameter θ that given x , have the max probability for making this y value happen: $\max p(y=1|x)$

最大似然估计, 通俗理解来说, 就是已经有样本信息 (数据), 得到最有可能导致这些样本结果出现的模型参数值!

$y^{\wedge}=p(y|x)$: for given x , the probability of output is $y=1$; since $y=1$ or $0 \rightarrow$
 $1-y^{\wedge} = p(y|x)$: for given x , the probability of output is $y=0$

>2. **output probability**: for $y=0$ and $y=1$, given signal input x ,

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

>3.: $\log (p(y|x)) = -J$

-->minimize J --> maximize $\log p, p$

(note: use log for numerical underflow problem)

>4. **Overall cost on m example**: assumption: all examples distribute independently.

$\log p(\cdot) = -J$.

$\log p$: use maximum likelihood estimate, find model parameter θ

J : minimize

by minimize cost function, is actually carrying out maximum likelihood with the logistic regression model, with assumption each training examples are identically independently distributed.

Cost on m examples

$$\log p(\text{labels in training set}) = \log \prod_{i=1}^m p(y^{(i)}|x^{(i)}) \leftarrow$$

$$\log p(\cdot) = \sum_{i=1}^m \log p(y^{(i)}|x^{(i)})$$

$$= \sum_{i=1}^m \left[-\sum_{j=1}^n \lambda_j y^{(i)} x_j^{(i)} \right]$$

$$\text{Cost: } J(w,b) = \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^n \lambda_j y^{(i)} x_j^{(i)}$$

Logistic regression cost function

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

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

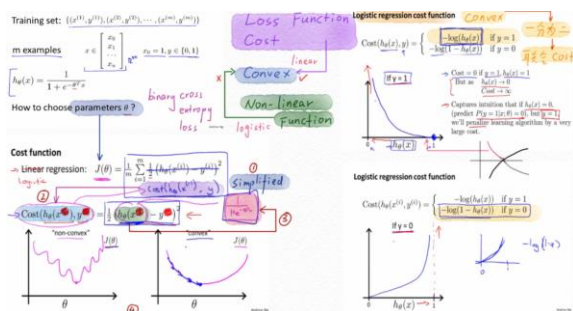
$$\text{If } y=1: p(y|x) = \hat{y} (1-\hat{y})^0$$

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

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

$$= -\sum_{j=1}^n \lambda_j \hat{y}^y x_j^{(i)}$$

Andrew



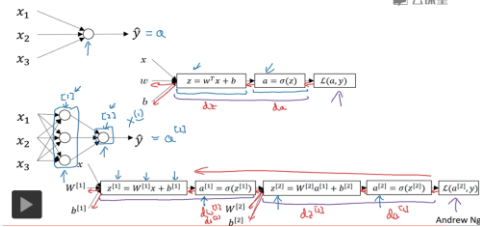
1. Neural network & logistic regression:

- >1. Logistic regression: just one node of neural network; one node in each layer
- >2. Neural network: logistic nodes stack together to become one layer, --> multiple layer --> neural network.

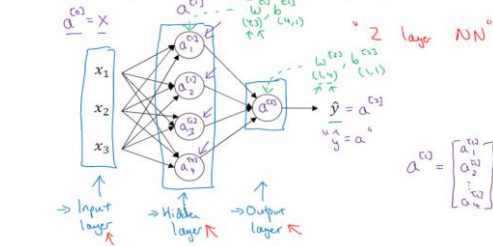
Neural network forward propagation and backpropagation: similar with logistic regression.

hidden layer: values here could not be observed; and have parameters associated.
Layer number: input layer not involved in accounting.

What is a Neural Network?



Neural Network Representation



2. Neural network computation detail_multi nodes, single training example:

- >1. every node in each layer: have two step computation:

$$Z = WT^* X + b$$

$$a = g(z)$$

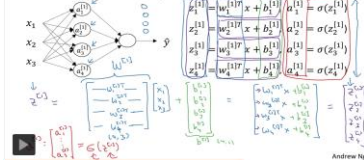
- >2. Vectorization for each layer
 $Z[i] = \text{layer } i \text{ all units stack together vertically}$
 $A[i] = g(Z[i])$

$$Z[i] = W[i] T^* X + b[i]$$

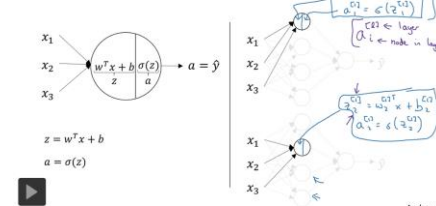
$$A[i] = g(Z[i])$$

Note: $b[i]$: $(n_i \times 1)$, have one value for each node in layer; $b[i]_i$: is real value, may different value for different nodes in the same layer.

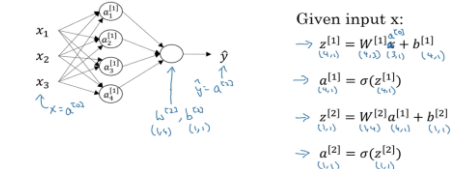
Neural Network Representation



Neural Network Representation



Neural Network Representation learning



3. Neural network on m training example_vectorization_multi nodes, multi example

similar with logistic regression:
 X : stacks in column for different example;

Z, A matrix:

- >1. stacks in column for different example

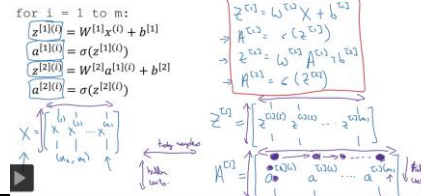
- >2. Matrix row/ horizontal direction: go through different examples;

- >3. matrix column/ vertical direction: go through different nodes in one layer.

Note:

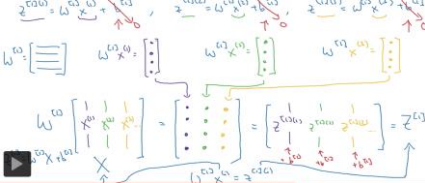
$b[i]$: broadcasting columnly (total m column = example number): same value for one node in layer, broadcasting in columns when have multiple examples.

Vectorizing across multiple examples

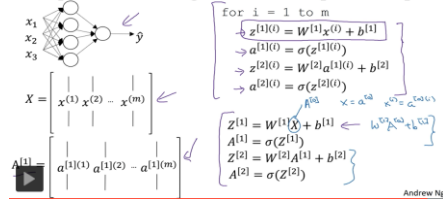


4. Neural network vectorization justification_multi node, multi example_two layer

Justification for vectorized implementation



Recap of vectorizing across multiple examples



5. Activation function: tanh

- >1. sigmoid function: 0-1;
- >2. tanh function: -1 ~1: shift function of sigmoid: go through (0,0) and then rescale to -1-1.

tanh works better than sigmoid, as have 0 mean / normalization effect: just incase might need normalize the data/weight???

-->

Hidden layer: use tanh function more popular than sigmoid

Output layer: will still use sigmoid for classification: output probability p: 0-1. while tanh is range from -1 to 1.

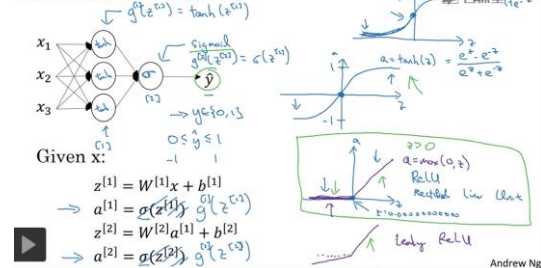
Note:

- >>1. activation function used in different (hidden) layer, maybe different

- >2. disadvantage of sigmoid, and tanh function: when z is large, gradient is close to 0, ---> gradient descent is very slow.---

(note: use batch norm could help improving sigmoid, tanh disadvantage: z large, gradient close to 0; batch norm --> $z \approx 0$)

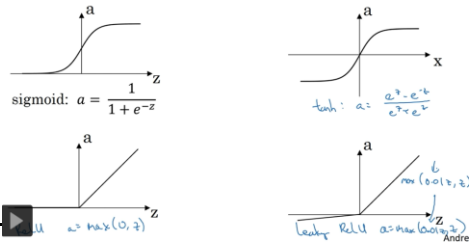
Activation functions



>3. Relu function: $\max(0, z)$
 no gradient when $z=0$, ok as in computation the chance of $z=0$ is too small;
 or could also define gradient value when $z=0$.
 disadvantage: derivative is equal to 0 when $z < 0$, in practice still work fine (have enough units to make $z > 0$), but have another version Relu: leaky Relu
 >4. leaky Relu: have slight slope when $z < 0$: this slope could be a parameter to learn, based on algorithm performance.
 works better than Relu, yet not used so much in practice

Activation function choosing rule:
 outlayer: if output is 0 or 1/binas classification, use sigmoid function.
 all other units: default choice is Relu: rectified linear unite

Pros and cons of activation functions



6. Why need non-linear activation function

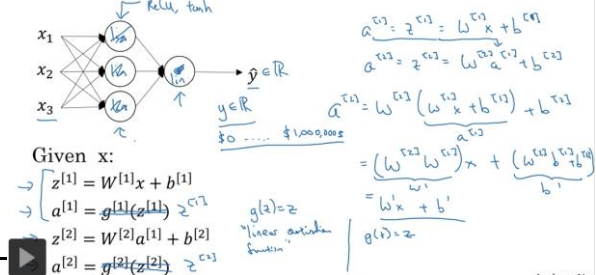
>1. If activation function $g[]$ is linear or no activation function in hidden layer & output layer -->

neural network is linear function no matter how many layers used/how many deeper, could not learn interesting property. In this case could remove hidden layers.

>2. if hidden layers activation function is linear, output layer is non-linear --> hidden layers could also removed as now hidden layer is more or less useless: two linear function combination of: input and hidden layer linear function, still linear.]

> just one place might use linear activation function: doing machine learning on a regression problem.
 e.g house price predict: output layer may use line activation function, but rarely happen to hidden layer. in this case output layer could also use non-linear activation function: Relu (price > 0).

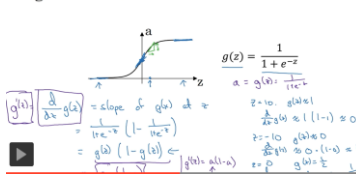
Activation function



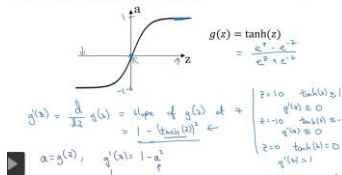
7. Derivation of activation function

Relu: derivation at $z=0$, ok to be any value set manually or could not define
 $y = \text{sigmoid} \rightarrow y' = y(1-y)$
 $y = \text{tanh} \rightarrow y' = 1 - y^2$

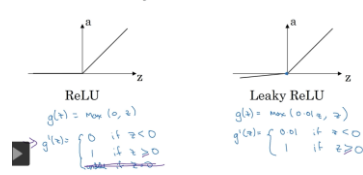
Sigmoid activation function



Tanh activation function



ReLU and Leaky ReLU



8. Neural network gradient descent

similar to logistic regression, cost function depending on activation function used.

>1. Forward propagation:

$Z[l] = W[l]^T A[l-1] + b[l]$
 $A[l] = g(Z[l])$

>2. Back propagation: on multi examples

$dA[l] = W[l+1]^T \cdot dZ[l+1]$
 $dZ[l] = dA[l] \cdot g'(Z[l])$

$dZ[l] = W[l+1]^T \cdot dZ[l+1] \cdot g'(Z[l])$ - combine above two step together

$W[l+1]^T \cdot dZ[l+1]$ is $(n, 1, m)$

$g'(Z[l])$ also is $(n, 1, m)$. * here means element wise product of two matrix.

$dW[l] = 1/m \cdot dZ[l] \cdot A[l-1]$

$db[l] = 1/m \cdot np.sum(dZ[l], axis=1, keepdims=True)$

Note:

1. keepdims=True: incase outputting funny rank 1 arrays, db[l] dimension is $(n, 1, 1)$

or could use reshape to make sure db[l] dimension.

2. In backpropagation: one tech. tip is to check matrix dimension mapping in computation.

Summary of gradient descent

$$dz^{[2]} = a^{[2]} - y$$

$$dW^{[2]} = dz^{[2]} a^{[1]T}$$

$$db^{[2]} = dz^{[2]}$$

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

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

$$db^{[1]} = dz^{[1]}$$

$$dZ^{[2]} = A^{[2]} - Y$$

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

$$db^{[2]} = \frac{1}{m} np.sum(dZ^{[2]}, axis=1, keepdims=True)$$

$$dZ^{[1]} = W^{[2]T} dZ^{[2]} * g^{[1]'}(Z^{[1]})$$

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} X^T$$

$$db^{[1]} = \frac{1}{m} np.sum(dZ^{[1]}, axis=1, keepdims=True)$$

9. Initialization

method1: Initialize weight W initialization: all to 0

>1. ok for logistic regression (as only have total one node)

>2. not work for neural network: hidden units in each layer is symmetry.

--> all units/nodes in each layer will compute same function: same as only have one node in each layer: all rows in layer weight $W[l]$, $W[l], i$ are same.

method 2: random initialize weight:

$W[l] = np.random.randn((2,2)) * 0.01$: randn-Gaussian distribution, range is 0-1;

factor 0.01 make sure $W[l]$ is very small-close to 0: because if initial weight is large, then $Z[l]$ may also very large-->when use tanh/sigmoid activation function, initial weight will be at the flat parts (activation function $\approx 1/-1$) of activation function, where $g'(Z)$ also to 0 -->will slow backpropagation, slow gradient descent, slow learning.

(note: batch norm will solve this problem-->BN make weight initialize not so important)

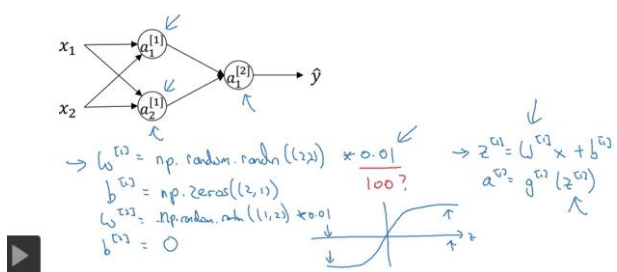
method3: b[l] could initialize all to 0: $b[l] = np.zeros((2,1))$

will not cause hidden units in each layer symmetry problem, which only caused when weight is set to 0.

Summary:

if initial weight is large, Z is large, will cause activation function tanh/sigmoid to be saturated, slowing down learning. But if have no tanh/sigmoid function, should not be a problem.

Random initialization



what happens if initialize only some layers weight to 0, not all layers $W[l]=0$. should not cause symmetric issue in layer units?

e.g:
two layer: layer 1: 2 nodes; layer 2: one node; training on one example

X-->layer 1 (2 nodes)---->layer 2 (one nodes)-->y^
note: layer 2 : like logistic regression.

Forward:

$Z[1] = W[1]T \cdot X + b[1];$
 $A[1] = g1(Z[1])$

$Z[2] = W[2]T \cdot A[1] + b[2];$
 $y^{\wedge} = A[2] = g2(Z[2])$

Backward:

$dA[2] <- -loss\ function$
 $dZ[2] = dA[2] \cdot g2'$

$dA[1] = W[2] \cdot dA[2] \cdot g2';$
 $dW[2] = 1/m \cdot dZ[2] \cdot A[1] = 1/m \cdot dA[2] \cdot g2' \cdot A[1]$
 $db[2] = 1/m \cdot np.sum(dZ[2], dim=1) = 1/m \cdot np.sum(dA[2] \cdot g2', dim=1)$

$dZ[2] = dA[2] \cdot g2'$

$dZ[1] = dA[1] \cdot g1' = W[2] \cdot dA[2] \cdot g2' \cdot g1';$
 $dW[1] = 1/m \cdot dZ[1] \cdot X = 1/m \cdot W[2] \cdot dA[2] \cdot g2' \cdot g1' \cdot X$
 $db[1] = 1/m \cdot np.sum(dZ[1], dim=1) = 1/m \cdot np.sum(W[2] \cdot dA[2] \cdot g2' \cdot g1', dim=1)$

E.G: $W1=W2= b1 = b2 =0$
forward:

$Z[1] = [0; 0]$
 $A[1] = [a;a] (a !=0)$

$Z[2] = W[2]T \cdot A[1] + b[2] = [0]$
 $y^{\wedge} = A[2] = g2(Z[2]) = [a2] != 0$

backprop;

$dA[2] = k(1 \times 1) <- -loss\ function$

-->**Second layer:**

$dZ[2] = dA[2] \cdot g2' = [k];$
 $dW[2] = 1/m \cdot dZ[2] \cdot A[1] = 1/m \cdot dA[2] \cdot g2' \cdot A[1] = 1/m \cdot k \cdot g2 \cdot [a;a]T = [dw2, dw2] != [0,0]$
 $db[2] = 1/m \cdot np.sum(dZ[2], dim=1) = 1/m \cdot np.sum(dA[2] \cdot g2', dim=1) = [db2] != 0$

-->second layer weight updated; yet symmetric $dW[2]$ is same for all units.

note: $dw[2]$: parameter same before and after updating;

-->**first layer:**

$dA[1] = W[2] \cdot dA[2] \cdot g2'; = [0;0]$
 $dZ[1] = dA[1] \cdot g1' = W[2] \cdot dA[2] \cdot g2' \cdot g1'; = [0;0]$
 $dW[1] = 1/m \cdot dZ[1] \cdot X = 1/m \cdot W[2] \cdot dA[2] \cdot g2' \cdot g1' \cdot X = 0$
 $db[1] = 1/m \cdot np.sum(dZ[1], dim=1) = 1/m \cdot np.sum(W[2] \cdot dA[2] \cdot g2' \cdot g1', dim=1) = 0$ (Note: db also influence by $W[i+1]$, yet ok to initialize b to 0.)

-->first layer:

no update as $dw1, db1=0$
error transfer from second layer to first layer by $W[2]$, while $W[2]$ initialize as all 0,-->error could not transfer to previous layer, previous layer parameter could not be updated.

1. Deep neural network

deep : have many hidden layers

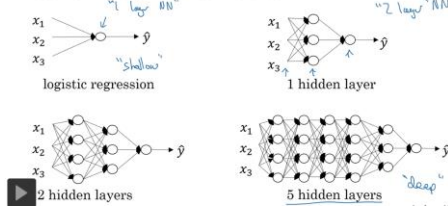
as hard to predict in advance exactly how deep a neural network needed, better try from logistic regression, then 1 hidden layer, 2 hidden layers, etc. and evaluate on dev set.

L = layer number: hidden layer + output layer

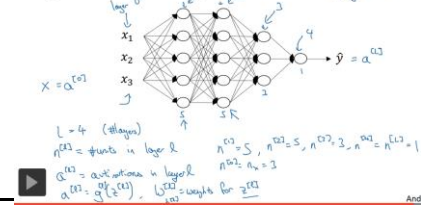
$n[l]$ = units number in layer l.

$a[l]$: activation function $a[0]=x$, $a[L]=y^{\wedge}$

What is a deep neural network?



Deep neural network notation



2. Deep neural network forward and back propagation on multi example same as shallow neural network

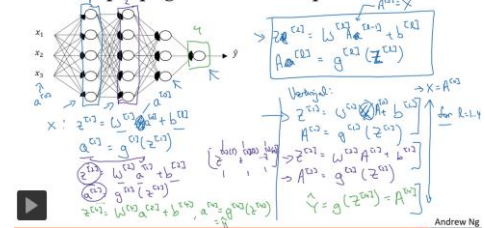
2.1. Forward propagation: in each layer: input $A[l-1]$, output $A[l]$;

$Z[l] = W[l] \cdot A[l-1] + b[l]$; (n, l, m); $A[0]=x$;

$A[l] = g(Z[l])$

need to repeat above for L layers, using 'for loop', which could not be vectorized.

Forward propagation in a deep network



2.2: forward and backward block: in each layer

>1. Forward propagation: input $A[l-1]$, output $A[l]$;

$Z[l] = W[l] \cdot A[l-1] + b[l]$;

$A[l] = g(Z[l])$

store cache $A[l]$, $Z[l]$, $W[l]$, $b[l]$ computed here, as will be used in backpropagation.

>2. Back propagation block: on multi examples: input $dA[l]$, output $dA[l-1]$

means give $a[l]$ change, how much will it affect $a[l-1]$

input: $dA[l]$

output:

> 1. $dZ[l] = dA[l] \cdot g'(Z[l])$

> 2. $dA[l-1] = W[l]^T \cdot dZ[l] = W[l]^T \cdot dA[l] \cdot g'(Z[l])$

>3. output: also

$dW[l] = 1/m \cdot dZ[l] \cdot A[l-1]$;

$db[l] = 1/m \cdot np.sum(dZ[l], axis=1, keepdims=True)$

>3. One iteration:

forward: $A[0] \rightarrow \dots \rightarrow A[L]$

Backward: $y^{\wedge} = A[L] \rightarrow \dots \rightarrow dA[L] \rightarrow \dots \rightarrow dW[1]$

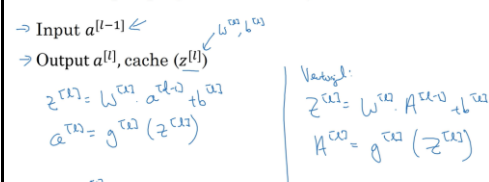
update W, b using gradient descent one time.

Note:

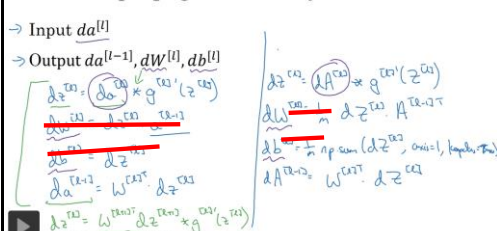
1. store each $Z[l]$, $W[l]$, $b[l]$ computed in forward, as will be used in backpropagation.

2. Forward, initialize by input x; while backward initialize by $dA[L]$, derivative of cost function over $A[L]$: (n, l, m)

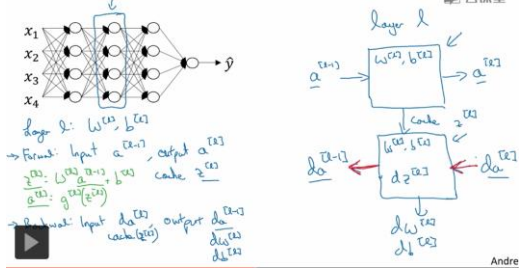
Forward propagation for layer l



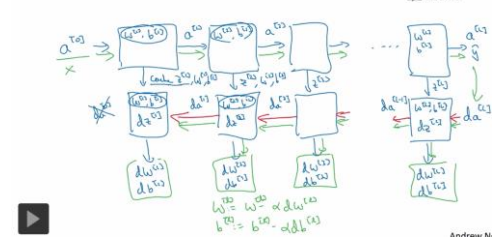
Backward propagation for layer l



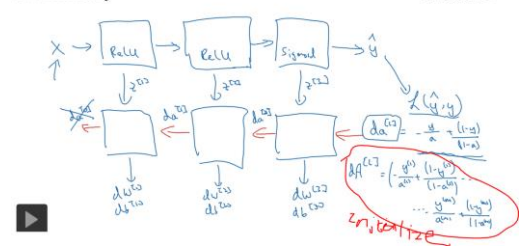
Forward and backward functions



Forward and backward functions



Summary



sum, or the average of all training set???

no matter loss function is the

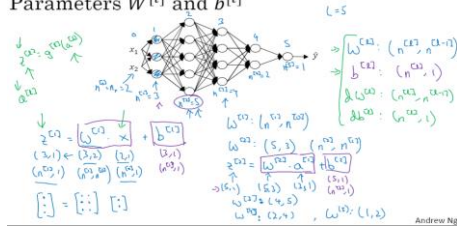
3. Matrix dimension check in neural_network_signal example

$a[l-1]$: input from layer $l-1$, dimension is $(n_l-1, 1)$
 $a[l]$: have n_l output, $(n_l, 1)$
 $z[l]$: $(n_l, 1)$: number of units in layer l ; same dimension as $a[l]$

$b[l]$: $(n_l, 1)$: number of units in layer l ;
 $W[l]$: $(n_l, n_{l-1}-1)$

$dw/db/dz/da$ have same dimension as W, b, Z, a
 $a[l]$ have same dimension with hidden status $z[l]$ (apply activation function on all hidden unit)

Parameters $W^{[l]}$ and $b^{[l]}$



4. Matrix dimension check in neural_network_multi example

$W[l]$ dimension no change: (n_l, n_{l-1}) ; dimension no influence by example size m
 $b[l]$ dimension no change: $(n_l, 1)$, yet python broadcasting will use change to (n_l, m) in matrix computation.

others a, z , changed, stack m example computation result in columns

$A[l-1]$: input from layer $l-1$, dimension is (n_{l-1}, m)

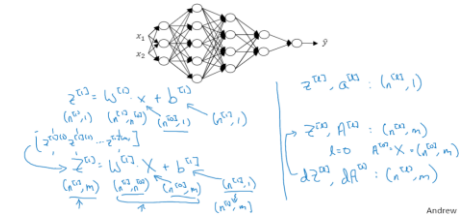
$A[l]$: have n_l output, (n_l, m)

$Z[l]$: (n_l, m) : number of hidden state in layer l ; same dimension as $A[l]$

$dw/db/dz/dA$ same dimension as W, b, Z, A

$A[l]$ have same dimension with $Z[l]$

Vectorized implementation



5. Deep neural network work well

5.1: intuition

e.g: face detection:

layer 1: may have 20 units to detect image edge.
 Layer 2: 40 units, group edges together to form parts of faces.
 Layer3: putting different face part to form face
 layer4: detect face

--->could assume earlier layers learning simple features, and later layer put learned feature together so neural network could learn on more complex functions.

eg: audio:

earlier layers may learn simple sound wave features, then later layer put together to form complex sound waveforms-->learn basic units of sound

5.2: circuit theory for deep learning

there are functions compute with a 'small' (small hidden units number in layers) L -layer deep neural network (more layers, but small units in each layer), that shall lower networks require exponentially more hidden units to computes (small layer, large hidden units in each layer).

e.g: compute XOR function: of x_1, \dots, x_n

-->use deep neuro network:

neural network depth is $\log(n)$, no much hidden units in each layer.

(note:

layer number $L: n/(2^L) = 1 \rightarrow L = \log(n) / \log(2)$

hidden units number in layer $l: n/(2^l)$, total hidden units number: $2n = 1+2+\dots+2^{L-1}$ ($2^L=n$)

)

>use 1 hidden layer neural network

hidden units is exponential large: $2^n(n-1)$ (have this number possible configurations/combinations).

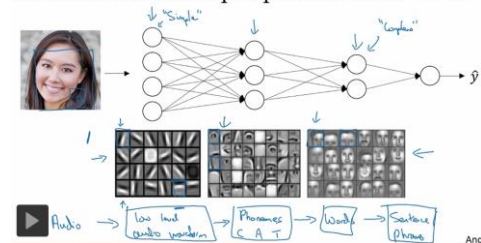
(note:

each feature have 2 value 0, 1-->input cases is: $2^n n$; hidden units could be $2^n(n-1)$

)

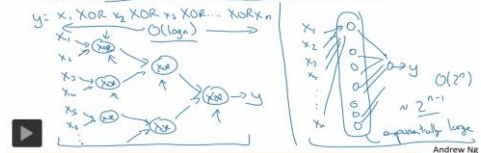
there are many mathematic functions that are much easier to compute with deep networks than with shallow networks.

Intuition about deep representation



Circuit theory and deep learning

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



6. Hyperparameters

parameter: W, b

6.1 Hyperparameter definition: control the ultimate parameter W and b that end up with.

- >1. learning rate α : determine how parameters evolve
- >2. Iteration number for gradient descent: batch size
- >3. Hidden layers number: L ,
- >4. hidden units n_l ,
- >5. activation function
- >6. momentum, minibatch size, regularizations
- etc.

6.2 Hyperparameter choose: based on empirical:

idea-->code-->experiment judge -->try idea...

evaluation different ideas based on cost function on dev set. many hyperparameters need to choose for new program, just try a range of values for them and see cost function/evaluation on dev set: based on performance: cost J/accuracy/signular evaluation F1 score-skewed negative and positive examples, etc.

(note: choice hyperameters to try: no grid choice-random sampling, coarse to fine and sampling 'uniformly' on good hyperparameter 'range'-rescaled is necessary -balance sentive to hyperamter change)

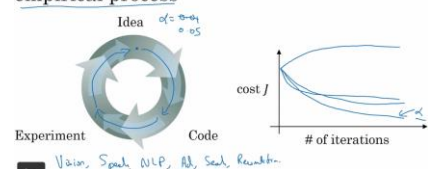
Note: best hyperparameter may change with time: computer infrastructure- CPU/GPU change. --->if working on a problem for an extended period of time for many years, then after a few months, try a few values for the hyper parameters and double check if there is a better value.

What are hyperparameters?

Parameters: $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, W^{[3]}, b^{[3]}, \dots$

Hyperparameter: learning rate α ,
 #iterations,
 #hidden layer L ,
 #hidden units n_1, n_2, \dots ,
 choice of activation function.
 Soft: Momentum, mini-batch size, regularizations.

Applied deep learning is a very empirical process



Forward and backward propagation

7. Neural network and human brain

no big relation between the two: may more easier to explain these equations's function when use 'brain'

And take logistic regression as a single neuron, is also very loose / simplistic analogy.

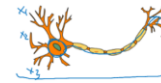
No one actually know what a single neuron does, even some of it's function similar to logistic regression.: no idea: how neuro learns, if similar with backpropagation, gradient descent.

$$\begin{aligned} Z^{[1]} &= W^{[1]}X + b^{[1]} \\ A^{[1]} &= g^{[1]}(Z^{[1]}) \\ Z^{[2]} &= W^{[2]}A^{[1]} + b^{[2]} \\ A^{[2]} &= g^{[2]}(Z^{[2]}) \\ &\vdots \\ A^{[L]} &= g^{[L]}(Z^{[L]}) = \hat{y} \end{aligned}$$

"It's like the brain."



$$\begin{aligned} dZ^{[L]} &= A^{[L]} - Y \\ dW^{[L]} &= \frac{1}{m} dZ^{[L]} A^{[L]T} \\ db^{[L]} &= \frac{1}{m} np.sum(dZ^{[L]}, axis = 1, keepdims = True) \\ dZ^{[L-1]} &= dW^{[L]T} dZ^{[L]} g'^{[L]}(Z^{[L-1]}) \\ dZ^{[1]} &= dW^{[L]T} dZ^{[2]} g'^{[1]}(Z^{[1]}) \\ dW^{[1]} &= \frac{1}{m} dZ^{[1]} A^{[1]T} \\ db^{[1]} &= \frac{1}{m} np.sum(dZ^{[1]}, axis = 1, keepdims = True) \end{aligned}$$



$x \rightarrow y$

AI

8. Neural network decision boundary??

-->decision boundary: only for classification problem, right?

-->classification problem: output use sigmoid or softmax function.

With learned paramter W, if use sigmoid activation function:

$y^{\wedge} = A[L] > c \rightarrow W[1] \cdot X = K$ (depending on c, W[2]...b) vector

-->decision boundary line number is row of W[1].

if overfitting (hidden units too many), --> hidden units > input feature: decision boundary line number: too many decision boundary line

when use other activation functions, decision boundary still be combination of different straight lines?-->depending on activation function, 'decision boundary' is more complex