# Neural Networks: Learning

## 1 Cost function

**Neural Network (Classification)**

$\rightarrow \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \ldots, (x^{(m)}, y^{(m)})\}$

$\rightarrow L =$ total no. of layers in network $\quad L = 4$

$\rightarrow s_l =$ no. of units (not counting bias unit) in layer $l \quad S_1 = 3, \quad S_2 = 5, \quad S_4 = S_L = 4$

Layer 1  Layer 2  Layer 3  Layer 4

**Binary classification**

$y = 0$ or $1 \leftarrow$

$h_\Theta(x)$

1 output unit $\leftarrow$

$h_\Theta(x) \in \mathbb{R}$

$S_L = 1. \quad K = 1 \leftarrow$

**Multi-class classification** (K classes)

$y \in \mathbb{R}^K$  E.g. $\begin{bmatrix}1\\0\\0\\0\end{bmatrix}, \begin{bmatrix}0\\1\\0\\0\end{bmatrix}, \begin{bmatrix}0\\0\\1\\0\end{bmatrix}, \begin{bmatrix}0\\0\\0\\1\end{bmatrix} \leftarrow$

pedestrian  car  motorcycle  truck

K output units

$h_\Theta(x) \in \mathbb{R}^k$

$S_l = K \qquad (k \geq 3)$

Andrew Ng

**Cost function**

Logistic regression:

$$J(\theta) = -\frac{1}{m}\left[\sum_{i=1}^{m} y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)}))\right] + \frac{\lambda}{2m}\sum_{j=1}^{n}\theta_j^2$$

Neural network:

$\rightarrow h_\Theta(x) \in \mathbb{R}^K \quad (h_\Theta(x))_i = i^{th}$ output

$$\rightarrow J(\Theta) = -\frac{1}{m}\left[\sum_{i=1}^{m}\sum_{k=1}^{K} y_k^{(i)} \log(h_\Theta(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_\Theta(x^{(i)}))_k)\right]$$

$$+ \frac{\lambda}{2m}\sum_{l=1}^{L-1}\sum_{i=1}^{s_l}\sum_{j=1}^{s_{l+1}}(\Theta_{ji}^{(l)})^2$$

$\Theta_{io}^{(1)} x_0 + \Theta_{ii}^{(1)} x_1 + \ldots$

$a_0$

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

$i = 0 \ldots S_l$

Andrew Ng

## 2 Backpropagation algorithm

**Gradient computation**

$$\rightarrow J(\Theta) = -\frac{1}{m}\left[\sum_{i=1}^{m}\sum_{k=1}^{K} y_k^{(i)} \log h_\theta(x^{(i)})_k + (1 - y_k^{(i)}) \log(1 - h_\theta(x^{(i)})_k)\right]$$

$$+ \frac{\lambda}{2m}\sum_{l=1}^{L-1}\sum_{i=1}^{s_l}\sum_{j=1}^{s_{l+1}}(\Theta_j^{(l)})^2$$

$\rightarrow \min_\Theta J(\Theta)$

Need code to compute:

$\rightarrow$ - $J(\Theta)$

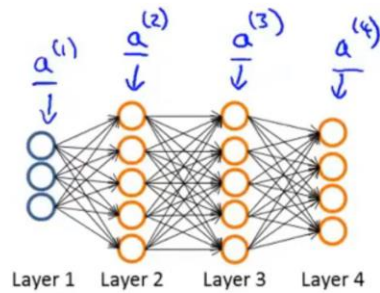$\rightarrow$ - $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$

$\Theta_{ij}^{(l)} \in \mathbb{R}$

## Gradient computation

Given one training example $(\underline{x, y})$:
Forward propagation:

$\quad a^{(1)} = \underline{x}$
$\rightarrow \underline{z^{(2)} = \Theta^{(1)} a^{(1)}}$
$\rightarrow a^{(2)} = g(z^{(2)}) \quad (\text{add } \underline{a_0^{(2)}})$
$\rightarrow z^{(3)} = \Theta^{(2)} a^{(2)}$
$\rightarrow a^{(3)} = g(z^{(3)}) \quad (\text{add } \underline{a_0^{(3)}})$
$\rightarrow z^{(4)} = \Theta^{(3)} a^{(3)}$
$\rightarrow a^{(4)} = h_\Theta(x) = g(z^{(4)})$

Layer 1    Layer 2    Layer 3    Layer 4

## Gradient computation: Backpropagation algorithm

Intuition: $\underline{\delta_j^{(l)}}$ = "error" of node $j$ in layer $l$.

$\delta^{(2)} \quad (l) \quad \delta^{(3)} \quad \delta^{(4)}$
$a_j$

Layer 1    Layer 2    Layer 3    Layer 4

For each output unit (layer L = 4)

$\boxed{\delta_j^{(4)} = \boxed{a_j^{(4)}} - y_j} \quad (h_\Theta(x))_j \quad \delta^{(4)} = a^{(4)} - y$

$\rightarrow \delta^{(3)} = (\Theta^{(3)})^T \delta^{(4)} .* g'(z^{(3)}) \qquad a^{(3)} .* (1 - a^{(3)})$

$\rightarrow \delta^{(2)} = (\Theta^{(2)})^T \delta^{(3)} .* g'(z^{(2)}) \qquad a^{(2)} .* (1 - a^{(2)})$

$(\text{No } \delta^{(1)})$

$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = a_j^{(l)} \delta_i^{(l+1)} \qquad (\text{ignoring } \lambda ; \text{ if } \lambda = 0) \leftarrow$

## Backpropagation algorithm

$\rightarrow$ Training set $\{(x^{(1)}, y^{(1)}), \ldots, (x^{(m)}, y^{(m)})\}$

Set $\underline{\triangle_{ij}^{(l)}} = 0$ (for all $l, i, j$). $\qquad (\text{used to compute } \frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta))$

For $\underline{i = 1 \text{ to } m} \leftarrow \quad (\underline{x^{(i)}, y^{(i)}})$.

$\quad$ Set $a^{(1)} = \underline{x^{(i)}}$

$\quad \rightarrow$ Perform forward propagation to compute $\underline{a^{(l)}}$ for $\underline{l = 2, 3, \ldots, L}$

$\quad \rightarrow$ Using $\underline{y^{(i)}}$, compute $\delta^{(L)} = \boxed{a^{(L)}} - \boxed{y^{(i)}}$

$\quad \rightarrow$ Compute $\underline{\delta^{(L-1)}, \delta^{(L-2)}, \ldots, \delta^{(2)}}$

$\quad \rightarrow \boxed{\triangle_{ij}^{(l)} := \triangle_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}} \leftarrow \qquad \triangle^{(l)} := \triangle^{(l)} + \delta^{(l+1)} (a^{(l)})^T.$

$\rightarrow \boxed{D_{ij}^{(l)}} := \frac{1}{m} \triangle_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)}$ if $j \neq 0 \qquad \frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = D_{ij}^{(l)}$

$\rightarrow \boxed{D_{ij}^{(l)}} := \frac{1}{m} \triangle_{ij}^{(l)} \qquad\qquad$ if $j = 0$

# 3 Backpropagation intuition

**Forward Propagation**



Nodes: $+1$, $x_1^{(i)}$, $x_2^{(i)}$ (input layer)

Layer 2: $z_1^{(2)} \rightarrow a_1^{(2)}$, $z_2^{(2)} \rightarrow a_2^{(2)}$, $+1$

Layer 3: $z_1^{(3)} \rightarrow a_1^{(3)}$, $z_2^{(3)} \rightarrow a_2^{(3)}$, $+1$

Layer 4: $z_1^{(4)} \rightarrow a_1^{(4)}$

$(x^{(i)}, y^{(i)})$

$$z_1^{(3)} = \Theta_{10}^{(2)} \times 1 + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} \cdot a_1^{(2)}$$

## What is backpropagation doing?

$$J(\Theta) = -\frac{1}{m}\left[\sum_{i=1}^{m} y^{(i)} \log(h_\Theta(x^{(i)})) + (1-y^{(i)})\log(1-(h_\Theta(x^{(i)})))\right]$$
$$+\frac{\lambda}{2m}\sum_{l=1}^{L-1}\sum_{i=1}^{s_l}\sum_{j=1}^{s_{l+1}}(\Theta_{ji}^{(l)})^2$$

$(x^{(i)}, y^{(i)})$

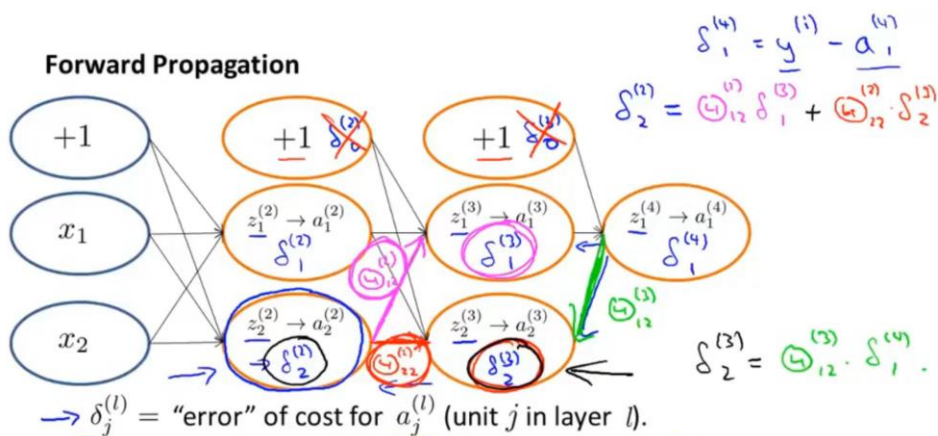Focusing on a single example $x^{(i)}$, $y^{(i)}$, the case of 1 output unit, and ignoring regularization ($\lambda = 0$),

$$\text{cost}(i) = y^{(i)} \log h_\Theta(x^{(i)}) + (1-y^{(i)})\log h_\Theta(x^{(i)})$$

(Think of $\text{cost}(i) \approx (h_\Theta(x^{(i)}) - y^{(i)})^2$)

I.e. how well is the network doing on example i?

**Forward Propagation**



$$\delta_1^{(4)} = y^{(i)} - a_1^{(4)}$$
$$\delta_2^{(2)} = \Theta_{12}^{(2)}\delta_1^{(3)} + \Theta_{22}^{(2)}\cdot\delta_2^{(3)}$$
$$\delta_2^{(3)} = \Theta_{12}^{(3)} \cdot \delta_1^{(4)}$$

Nodes layer 2: $z_1^{(2)} \rightarrow a_1^{(2)}$, $\delta_1^{(2)}$; $z_2^{(2)} \rightarrow a_2^{(2)}$, $\delta_2^{(2)}$

Layer 3: $z_1^{(3)} \rightarrow a_1^{(3)}$, $\delta_1^{(3)}$; $z_2^{(3)} \rightarrow a_2^{(3)}$, $\delta_2^{(3)}$

Layer 4: $z_1^{(4)} \rightarrow a_1^{(4)}$, $\delta_1^{(4)}$

$\delta_j^{(l)} = $ "error" of cost for $a_j^{(l)}$ (unit $j$ in layer $l$).

Formally, $\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(i)$ (for $j \geq 0$), where
$\text{cost}(i) = y^{(i)} \log h_\Theta(x^{(i)}) + (1-y^{(i)})\log h_\Theta(x^{(i)})$

# 4 Implementation note: Unrolling parameters

## Advanced optimization

```
function [jVal, gradient] = costFunction(theta)
    ...
```
$\hookrightarrow \mathbb{R}^{n+1}$  $\hookrightarrow \mathbb{R}^{n+1}$ (vectors)

```
optTheta = fminunc(@costFunction, initialTheta, options)
```

Neural Network (L=4):

$\rightarrow \Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ - matrices (`Theta1`, `Theta2`, `Theta3`)

$\rightarrow D^{(1)}, D^{(2)}, D^{(3)}$ - matrices (`D1`, `D2`, `D3`)

"Unroll" into vectors

### Example

$s_1 = 10, s_2 = 10, s_3 = 1$

$\rightarrow \Theta^{(1)} \in \mathbb{R}^{10 \times 11}, \Theta^{(2)} \in \mathbb{R}^{10 \times 11}, \Theta^{(3)} \in \mathbb{R}^{1 \times 11}$

$\rightarrow D^{(1)} \in \mathbb{R}^{10 \times 11}, D^{(2)} \in \mathbb{R}^{10 \times 11}, D^{(3)} \in \mathbb{R}^{1 \times 11}$
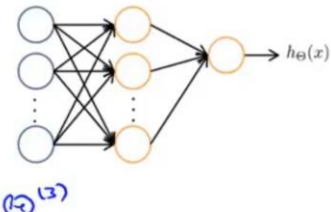
$\Theta^{(1)}$   $\Theta^{(2)}$   $\Theta^{(3)}$

```
thetaVec = [ Theta1(:); Theta2(:); Theta3(:)];
DVec = [D1(:); D2(:); D3(:)];

Theta1 = reshape(thetaVec(1:110),10,11);
Theta2 = reshape(thetaVec(111:220),10,11);
Theta3 = reshape(thetaVec(221:231),1,11);
```

```
Octave-3.2.4

  1  1  1  1  1  1  1  1  1  1  1
  1  1  1  1  1  1  1  1  1  1  1
  1  1  1  1  1  1  1  1  1  1  1
  1  1  1  1  1  1  1  1  1  1  1
  1  1  1  1  1  1  1  1  1  1  1
  1  1  1  1  1  1  1  1  1  1  1
  1  1  1  1  1  1  1  1  1  1  1

>> reshape(thetaVec(111:220), 10,11)
ans =

  2  2  2  2  2  2  2  2  2  2  2
  2  2  2  2  2  2  2  2  2  2  2
  2  2  2  2  2  2  2  2  2  2  2
  2  2  2  2  2  2  2  2  2  2  2
  2  2  2  2  2  2  2  2  2  2  2
  2  2  2  2  2  2  2  2  2  2  2
  2  2  2  2  2  2  2  2  2  2  2
  2  2  2  2  2  2  2  2  2  2  2
  2  2  2  2  2  2  2  2  2  2  2
  2  2  2  2  2  2  2  2  2  2  2

>> reshape(thetaVec(221:231), 1,11)
```
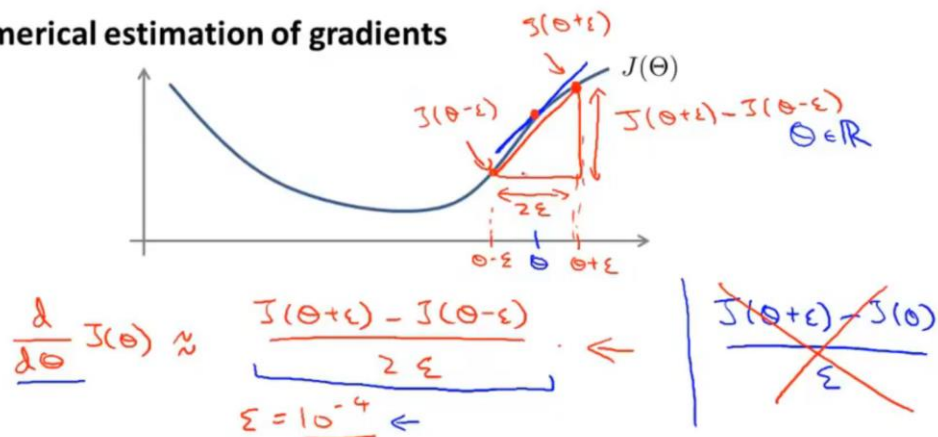5:45 PM
10/16/2011

**Learning Algorithm**

→ Have initial parameters $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$.

→ Unroll to get `initialTheta` to pass to

→ `fminunc(@costFunction, initialTheta, options)`

```
function [jval, gradientVec] = costFunction(thetaVec)
```
→ From `thetaVec`, get $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$.  reshape

→ Use forward prop/back prop to compute $D^{(1)}, D^{(2)}, D^{(3)}$ and $J(\Theta)$.
Unroll $D^{(1)}, D^{(2)}, D^{(3)}$ to get `gradientVec`.

# 5 Gradient checking

**Numerical estimation of gradients**



$$\frac{d}{d\Theta} J(\Theta) \approx \frac{J(\Theta+\varepsilon) - J(\Theta-\varepsilon)}{2\varepsilon} \leftarrow$$

$$\varepsilon = 10^{-4} \leftarrow$$

$$\frac{J(\Theta+\varepsilon) - J(\Theta)}{\varepsilon}$$

Implement: `gradApprox = (J(theta + EPSILON) - J(theta - EPSILON))`
`/ (2*EPSILON)`

**Parameter vector $\theta$**

→ $\theta \in \mathbb{R}^n$ (E.g. $\theta$ is "unrolled" version of $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$)

→ $\theta = [\theta_1, \theta_2, \theta_3, \ldots, \theta_n]$

→ $\frac{\partial}{\partial \theta_1} J(\theta) \approx \frac{J(\theta_1+\epsilon, \theta_2, \theta_3, \ldots, \theta_n) - J(\theta_1-\epsilon, \theta_2, \theta_3, \ldots, \theta_n)}{2\epsilon}$

→ $\frac{\partial}{\partial \theta_2} J(\theta) \approx \frac{J(\theta_1, \theta_2+\epsilon, \theta_3, \ldots, \theta_n) - J(\theta_1, \theta_2-\epsilon, \theta_3, \ldots, \theta_n)}{2\epsilon}$

$\vdots$

→ $\frac{\partial}{\partial \theta_n} J(\theta) \approx \frac{J(\theta_1, \theta_2, \theta_3, \ldots, \theta_n+\epsilon) - J(\theta_1, \theta_2, \theta_3, \ldots, \theta_n-\epsilon)}{2\epsilon}$

```
for i = 1:n,  ←
    thetaPlus = theta;
    thetaPlus(i) = thetaPlus(i) + EPSILON;
    thetaMinus = theta;
    thetaMinus(i) = thetaMinus(i) - EPSILON;
    gradApprox(i) = (J(thetaPlus) - J(thetaMinus))
                    /(2*EPSILON);
end;
```

$$\begin{bmatrix} \Theta_1 \\ \Theta_2 \\ \Theta_i + \varepsilon \\ \vdots \\ \Theta_n \end{bmatrix} \rightarrow \Theta_i - \varepsilon$$

$$\frac{\partial}{\partial \Theta_i} J(\Theta).$$

Check that $\boxed{\text{gradApprox}} \approx \boxed{\text{DVec}}$  ←

↑

From backprop.

**Implementation Note:**
- Implement backprop to compute DVec (unrolled $D^{(1)}, D^{(2)}, D^{(3)}$).
- Implement numerical gradient check to compute gradApprox.
- Make sure they give similar values.
- Turn off gradient checking. Using backprop code for learning.

$$\delta^{(4)}, \delta^{(3)}, \delta^{(2)} \rightarrow \text{DVec}$$

**Important:**
- Be sure to disable your gradient checking code before training your classifier. If you run numerical gradient computation on every iteration of gradient descent (or in the inner loop of costFunction(...))your code will be very slow.

# 6 Random initialization

**Initial value of** $\Theta$

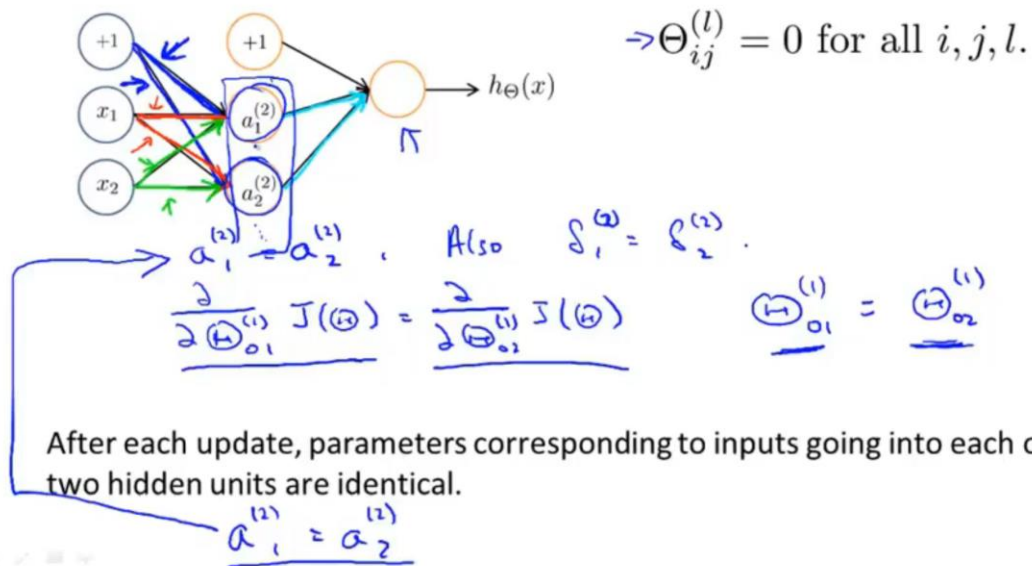For gradient descent and advanced optimization method, need initial value for $\Theta$.

```
optTheta = fminunc(@costFunction,
            initialTheta, options)
```

Consider gradient descent
Set `initialTheta = zeros(n,1)` ?

## Zero initialization



$$\Theta_{ij}^{(l)} = 0 \text{ for all } i, j, l.$$

$a_1^{(2)} = a_2^{(2)}$ , Also $\delta_1^{(2)} = \delta_2^{(2)}$ .

$$\frac{\partial}{\partial \Theta_{01}^{(i)}} J(\Theta) = \frac{\partial}{\partial \Theta_{02}^{(i)}} J(\Theta) \qquad \Theta_{01}^{(1)} = \Theta_{02}^{(1)}$$

After each update, parameters corresponding to inputs going into each of two hidden units are identical.

$$a_1^{(2)} = a_2^{(2)}$$

## Random initialization: Symmetry breaking

→ Initialize each $\Theta_{ij}^{(l)}$ to a random value in $[-\epsilon, \epsilon]$
(i.e. $-\epsilon \le \Theta_{ij}^{(l)} \le \epsilon$)

E.g.

random $10 \times 11$ matrix (betw. 0 and 1)

→ `Theta1 = rand(10,11)*(2*INIT_EPSILON)`
`- INIT_EPSILON;`    $[-\epsilon, \epsilon]$

→ `Theta2 = rand(1,11)*(2*INIT_EPSILON)`
`- INIT_EPSILON;`

# 7 Putting it together

### Training a neural network
Pick a network architecture (connectivity pattern between neurons)



→ No. of input units: Dimension of features $x^{(i)}$
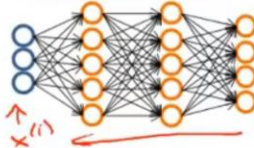→ No. output units: Number of classes
Reasonable default: 1 hidden layer, or if >1 hidden layer, have same no. of hidden units in every layer (usually the more the better)

$y \in \{1, 2, 3, \ldots, 10\}$

$$y = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \text{ or } \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

## Training a neural network

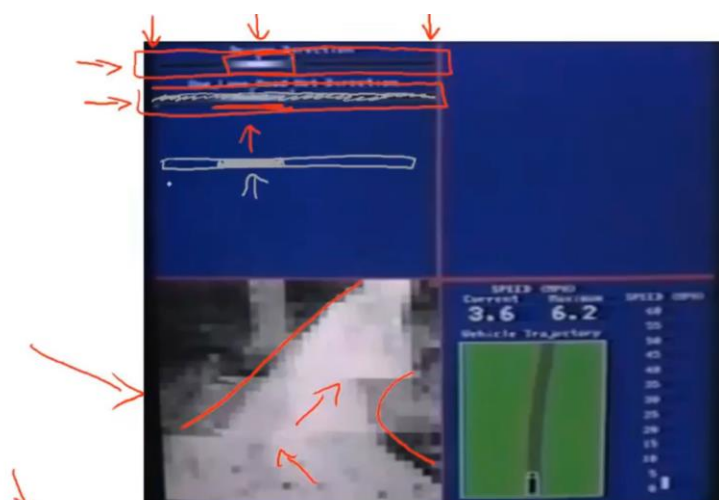→ 1. Randomly initialize weights
→ 2. Implement forward propagation to get $h_\Theta(x^{(i)})$ for any $x^{(i)}$
→ 3. Implement code to compute cost function $J(\Theta)$
→ 4. Implement backprop to compute partial derivatives $\dfrac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$

→ **for** i = 1:m {   $(x^{(1)}, y^{(1)})$   $(x^{(2)}, y^{(2)})$, ..... $(x^{(m)}, y^{(m)})$

→ Perform forward propagation and backpropagation using example $(x^{(i)}, y^{(i)})$

(Get activations $a^{(l)}$ and delta terms $\delta^{(l)}$ for $l = 2, \dots, L$).

$\Delta^{(2)} := \Delta^{(2)} + \delta^{(3)} (a^{(2)})^T$

$\vdots$

compute $\dfrac{\partial}{\partial \Theta_{jk}^{(2)}} J(\Theta)$.

$x^{(i)}$

An

## Training a neural network

→ 5. Use gradient checking to compare $\dfrac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$ computed using backpropagation vs. using numerical estimate of gradient of $J(\Theta)$.

→ Then disable gradient checking code.

→ 6. Use gradient descent or advanced optimization method with backpropagation to try to minimize $J(\Theta)$ as a function of parameters $\Theta$

$\dfrac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$

$J(\Theta)$   —   non-convex.

# 8 Autonomous driving example