

8. Multi-class Classification

Most of this material is from Prof. Andrew Ng and Chang's slides.

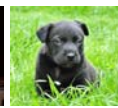
Recognizing cats, dogs, baby chicks, and others



3



1



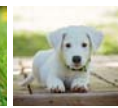
2



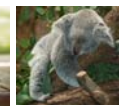
0



3



2



0

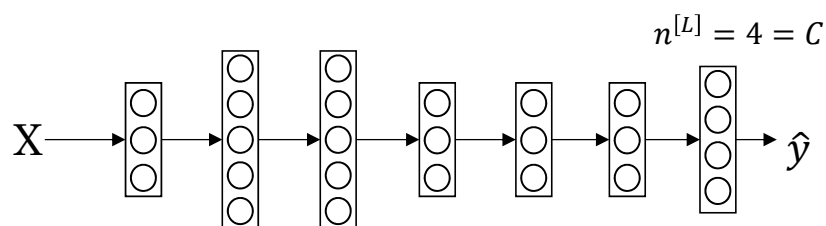


1

Recognizing cats, dogs, baby chicks, and others



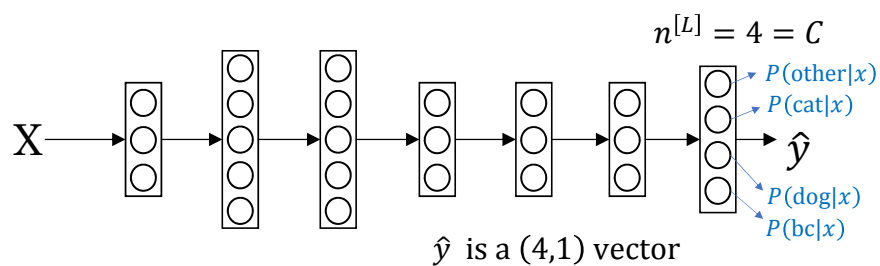
$\mathcal{C} = \text{\#classes} = 4$



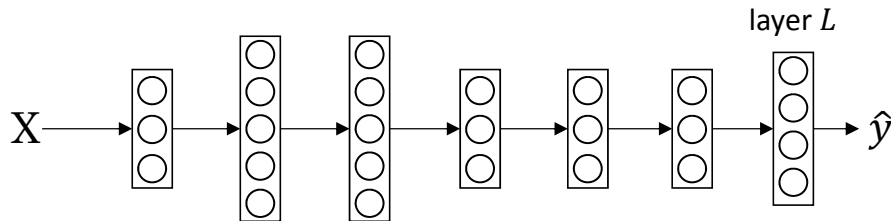
Recognizing cats, dogs, baby chicks, and others



$\mathcal{C} = \text{\#classes} = 4$



Recognizing cats, dogs, baby chicks, and others



(4,1)

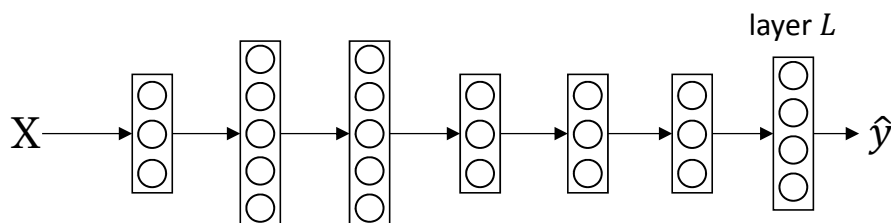
$$Z^{[L]} = W^{[L]}a^{[L-1]} + b^{[L]}$$

Activation function:

(4,1) $t = e^{(Z^{[L]})}$

(4,1) $a^{[L]} = \frac{e^{(Z^{[L]})}}{\sum_{j=1}^4 t_j} \rightarrow a_i^{[L]} = \frac{e^{(z_i^{[L]})}}{\sum_{j=1}^4 e^{(z_j^{[L]})}}$

Recognizing cats, dogs, baby chicks, and others



(4,1)

$$Z^{[L]} = W^{[L]}a^{[L-1]} + b^{[L]}$$

Activation function:

(4,1) $t = e^{(Z^{[L]})}$

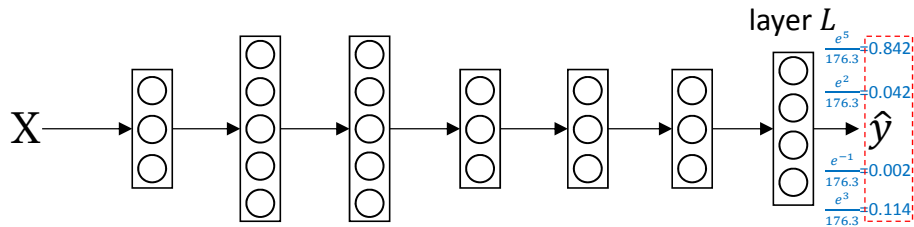
(4,1) $a^{[L]} = \frac{e^{(Z^{[L]})}}{\sum_{j=1}^4 t_j} \rightarrow a_i^{[L]} = \frac{e^{(z_i^{[L]})}}{\sum_{j=1}^4 e^{(z_j^{[L]})}}$

$$Z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}$$

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

$$a^{[L]} = \frac{t}{176.3}$$

Recognizing cats, dogs, baby chicks, and others



(4,1)

$$Z^{[L]} = W^{[L]}a^{[L-1]} + b^{[L]}$$

Activation function:

(4,1) $t = e^{(Z^{[L]})}$

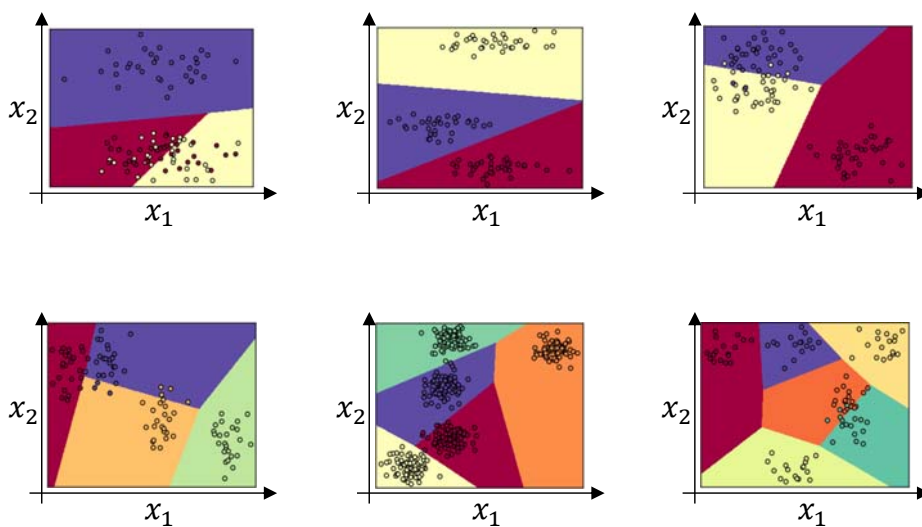
(4,1) $a^{[L]} = \frac{e^{(Z^{[L]})}}{\sum_{j=1}^4 t_j} \rightarrow a_i^{[L]} = \frac{e^{(z_i^{[L]})}}{\sum_{j=1}^4 e^{(z_j^{[L]})}}$
 $= \hat{y}$

$$Z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}$$

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

$$a^{[L]} = \frac{t}{176.3}$$

(1-layer) Softmax classifier examples



Understanding softmax

$$z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix}$$

$$g^{[L]}(z^{[L]}) = \begin{bmatrix} e^5/(e^5 + e^2 + e^{-1} + e^3) \\ e^2/(e^5 + e^2 + e^{-1} + e^3) \\ e^{-1}/(e^5 + e^2 + e^{-1} + e^3) \\ e^3/(e^5 + e^2 + e^{-1} + e^3) \end{bmatrix} = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix}$$

Understanding softmax

$$z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix}$$

$$g^{[L]}(z^{[L]}) = \begin{bmatrix} e^5/(e^5 + e^2 + e^{-1} + e^3) \\ e^2/(e^5 + e^2 + e^{-1} + e^3) \\ e^{-1}/(e^5 + e^2 + e^{-1} + e^3) \\ e^3/(e^5 + e^2 + e^{-1} + e^3) \end{bmatrix} \stackrel{\text{"soft max"}}{=} \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix} \quad \text{vs.} \quad \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \stackrel{\text{"hard max"}}{=}$$

Understanding softmax

$$z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix}$$

$$g^{[L]}(z^{[L]}) = \begin{bmatrix} e^5/(e^5 + e^2 + e^{-1} + e^3) \\ e^2/(e^5 + e^2 + e^{-1} + e^3) \\ e^{-1}/(e^5 + e^2 + e^{-1} + e^3) \\ e^3/(e^5 + e^2 + e^{-1} + e^3) \end{bmatrix} \overset{\text{"soft max"}}{=} \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix} \quad \text{vs.} \quad \overset{\text{"hard max"}}{\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}}$$

Softmax regression generalizes logistic regression to C classes.

Understanding softmax

$$z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix}$$

$$g^{[L]}(z^{[L]}) = \begin{bmatrix} e^5/(e^5 + e^2 + e^{-1} + e^3) \\ e^2/(e^5 + e^2 + e^{-1} + e^3) \\ e^{-1}/(e^5 + e^2 + e^{-1} + e^3) \\ e^3/(e^5 + e^2 + e^{-1} + e^3) \end{bmatrix} \overset{\text{"soft max"}}{=} \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix} \quad \text{vs.} \quad \overset{\text{"hard max"}}{\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}}$$

Softmax regression generalizes logistic regression to C classes.

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

Loss function

assume

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad a^{[L]} = \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

Loss function

assume

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad a^{[L]} = \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

$$\mathcal{L}(\hat{y}, y) = - \sum_{j=1}^4 y_j \log \hat{y}_j$$

Loss function

assume

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad a^{[L]} = \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

$$\mathcal{L}(\hat{y}, y) = - \sum_{j=1}^4 y_j \log \hat{y}_j$$

$$= -y_2 \log \hat{y}_2 = -\log \hat{y}_2$$



to make $\mathcal{L}(\hat{y}, y)$ small,
we should make \hat{y}_2 big.

minimizing loss function is **justified by maximum likelihood principle**

→ find $\operatorname{argmin}_{w,b} (-\log \hat{y}_i) = \operatorname{argmax}_{w,b} (\hat{y}_i) = \operatorname{argmax}_{w,b} (p(y_i|W, b))$ such that $y_i = 1$

Loss function

assume

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad a^{[L]} = \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

$$\mathcal{L}(\hat{y}, y) = - \sum_{j=1}^4 y_j \log \hat{y}_j$$

$$= -y_2 \log \hat{y}_2 = -\log \hat{y}_2$$



to make $\mathcal{L}(\hat{y}, y)$ small,
we should make \hat{y}_2 big.

$$J(W^{[1]}, b^{[1]}, \dots) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

Loss function

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad a^{[L]} = \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

$$\begin{aligned} \mathcal{L}(\hat{y}, y) &= - \sum_{j=1}^4 y_j \log \hat{y}_j \\ &= -y_2 \log \hat{y}_2 = -\log \hat{y}_2 \end{aligned}$$



to make $\mathcal{L}(\hat{y}, y)$ small,
we should make \hat{y}_2 big.

$$J(W^{[1]}, b^{[1]}, \dots) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

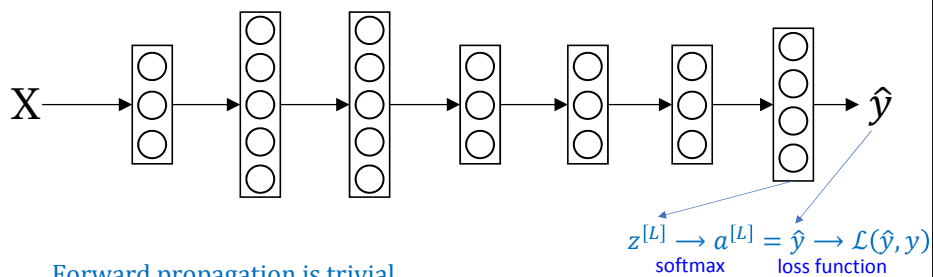
$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(m)} \end{bmatrix}$$

(4, m)

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

(4, m)

Gradient descent with softmax



Forward propagation is trivial.

Backpropagation: $\frac{dz^{[L]}}{\partial z^{[L]}} = \hat{y} - y$

(4,1)

Proof of $dz^{[L]} = \frac{\partial J}{\partial z^{[L]}} = \hat{y} - y$ $\quad da = \frac{\partial J}{\partial a} = p - y$

For simplicity, set $p = \hat{y}$ and $p_i = \hat{y}_i$

$a = z^{[L]}$ and $a_i = z_i^{[L]}$

$\therefore p_i = \frac{e^{(a_i)}}{\sum_k e^{(a_k)}}$

$$\begin{aligned} \frac{\partial p_i}{\partial a_i} &= \frac{\partial \frac{\exp(a_i)}{\sum_k \exp(a_k)}}{\partial a_i} \\ &= \frac{\exp(a_i) \sum_k \exp(a_k) - \exp(a_i) \exp(a_i)}{(\sum_k \exp(a_k))^2} \\ &= \frac{\exp(a_i) [\sum_k \{\exp(a_k)\} - \exp(a_i)]}{(\sum_k \exp(a_k))^2} \\ &= \frac{\exp(a_i)}{\sum_k \exp(a_k)} \frac{\sum_k \{\exp(a_k)\} - \exp(a_i)}{\sum_k \exp(a_k)} \\ &= \frac{\exp(a_i)}{\sum_k \exp(a_k)} \left(1 - \frac{\exp(a_i)}{\sum_k \exp(a_k)}\right) \\ &= p_i (1 - p_i) \end{aligned}$$

for $i \neq j$, $\frac{\partial p_i}{\partial a_j} = \frac{0 - \exp(a_i) \exp(a_j)}{(\sum_k \exp(a_k))^2} = -\frac{\exp(a_i)}{\sum_k \exp(a_k)} \frac{\exp(a_j)}{\sum_k \exp(a_k)} = -p_i p_j$

Proof of $da = \frac{\partial J}{\partial a} = p - y$

$$\begin{aligned} \frac{\partial J}{\partial a_i} &= \frac{\partial (-\sum_j y_j \log p_j)}{\partial a_i} \\ &= -\sum_j y_j \frac{\partial \log p_j}{\partial a_i} \\ &= -\sum_j y_j \frac{1}{p_j} \frac{\partial p_j}{\partial a_i} = -\frac{y_i}{p_i} (1 - p_i) - \sum_{i \neq j} \frac{y_j}{p_j} (-p_i p_j) \\ &= -y_i + y_i p_i + \sum_{i \neq j} y_j p_i \\ &= -y_i + \sum_j y_j p_i \\ &= -y_i + p_i \sum_j y_j \\ &= p_i - y_i \\ \therefore \frac{\partial J}{\partial a} &= [p_i - y_i] = p - y \implies dz^{[L]} = \frac{\partial J}{\partial z^{[L]}} = \hat{y} - y \end{aligned}$$