

Deep Learning

Supervised learning w/ non-linear models

before $h_{\theta}(x) = \theta^T \phi(x)$ non-linear in x
Kernel method linear in θ

$$\text{eg. } h_{\theta}(x) = \sqrt{\theta_1^3 + \theta_2 x_3} + \sqrt{\theta_4} x_7$$

dataset $\{(x^{(i)}, y^{(i)})\}_{i=1}^n, x^{(i)} \in \mathbb{R}^d, y^{(i)} \in \mathbb{R}$
 $h_{\theta}(x) : \mathbb{R}^d \rightarrow \mathbb{R}$

cost/loss f^n

$$J^{(i)}(\theta) = (y^{(i)} - h_{\theta}(x^{(i)}))^2 \quad \text{squared loss}$$

cost f^n for entire dataset

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n J^{(i)}(\theta)$$

optimization objective

$$\min_{\theta} J(\theta)$$

gradient descent

$$\theta := \theta - \alpha \nabla J(\theta)$$

Stochastic Gradient Descent (SGD)

for $i = 1$ to N_{iter}

Sample j from $\{1, \dots, n\}$ randomly

$$\theta := \theta - \alpha \nabla J^{(j)}(\theta)$$

minibatch SGD

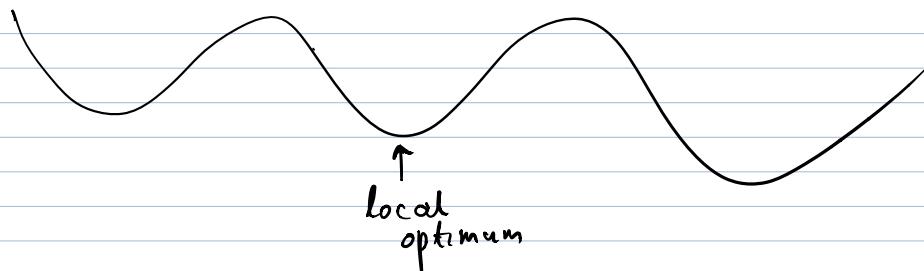
Idea: Computing B gradients $\nabla J^{(j_1)}(\theta), \dots, \nabla J^{(j_B)}(\theta)$

together is faster than sequential computations

for $i = 1$ to n_{iter}

Sample B examples $\{j_1 \dots j_B\}$ from $\{1 \dots n\}$

$$\theta := \theta - \frac{\alpha}{B} \sum_{k=1}^B \nabla J^{(j_k)}(\theta)$$



remaining questions

① How to define $h_\theta(x)$?

② how to compute gradient

— Logistic Regression

— Neural Networks

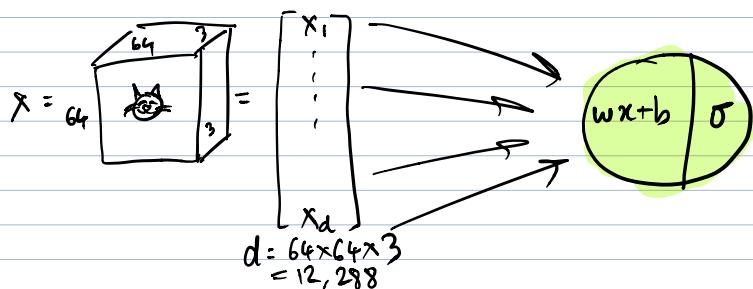
— computational power

— data available

— algorithms

① Logistic Regression

goal Find cats in images $1 \rightarrow$ presence of cat
 $0 \rightarrow$ absence of cat



$$\begin{aligned}\hat{y} &= \sigma(w^T x) \\ &= \sigma(wx + b) \\ &\downarrow \\ &1 \times 12288\end{aligned}$$

(i) initialize w , b

(ii) Find optimal w , b

(iii) Use $\hat{y} = \sigma(wx+b)$ to predict

$$\lambda = -[y \log \hat{y} + (1-y) \log (1-\hat{y})]$$

$$w = w - \alpha \frac{\partial \lambda}{\partial w}$$

$$b = b - \alpha \frac{\partial \lambda}{\partial b}$$

$$\# \text{parameters} = d+1$$

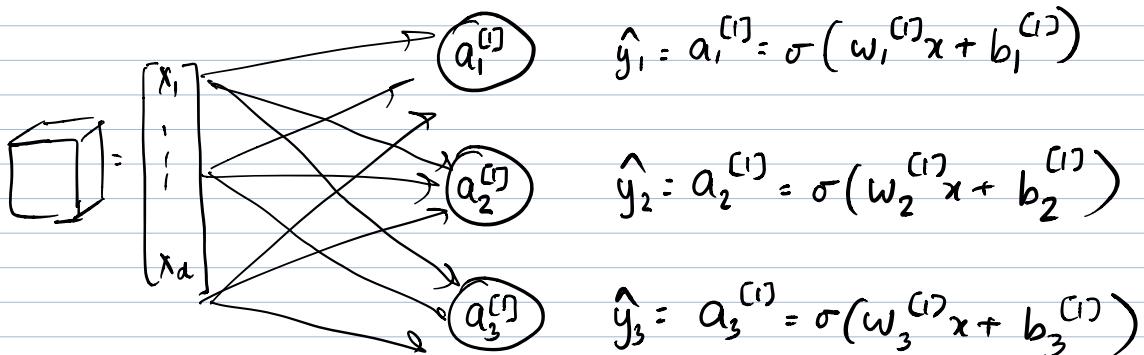
Note: fixed typo in the live notes

$$12288+1$$

neuron = linear + activation

model = architecture + parameters

Goal 2.0: Find cat / lion / iguana in images



square brackets [1] \equiv layer

subscripts \equiv identify neuron within layer

$a_1^{(1)}$ \leftarrow 1st layer

$a_2^{(1)}$ \leftarrow 2nd neuron in 1st layer

$$\# \text{ param} = 3(d+1)$$

images + labels



$$= \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \begin{array}{l} \text{--- cat} \\ \text{--- lion} \\ \text{--- iguana} \end{array}$$



$$= \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \begin{array}{l} \text{cat} \\ \text{lion} \\ \text{iguana} \end{array}$$

Goal 3.0 :

add constraint : unique animal in image



$$= \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

$$z_1^{(0)}$$

$$\hat{y}_1 = e^{z_1^{(0)}} / \sum_{k=1}^3 e^{z_k^{(0)}}$$

$$z_2^{(0)}$$

$$\hat{y}_2 = e^{z_2^{(0)}} / \sum_{k=1}^3 e^{z_k^{(0)}}$$

$$z_3^{(0)}$$

$$\hat{y}_3 = e^{z_3^{(0)}} / \sum_{k=1}^3 e^{z_k^{(0)}}$$

Note: Fixed typo in live notes

$$\hat{y}_1 = \sigma(w_1^{(0)} x + b_1^{(0)})$$

$\underbrace{w_1^{(0)} x + b_1^{(0)}}_{z_1^{(0)}}$

softmax
multiclass
network

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

cat

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

lion

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

iguana

$$\# \text{ parameters} = 3(d+1)$$

Loss function

$$L_{3N} = - \sum_{k=1}^3 [y_k \log \hat{y}_k + (1-y_k) \log (1-\hat{y}_k)]$$

Binary
Cross Entropy
Loss Function

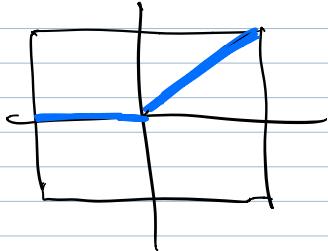
$$L_{CE} = - \sum_{k=1}^3 y_k \log \hat{y}_k$$

Q^n: Predict age of cat instead of presence of cat?

Options

① Several neurons to predict different ages

② Change activation fn



$$f(x) = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases}$$

ReLU

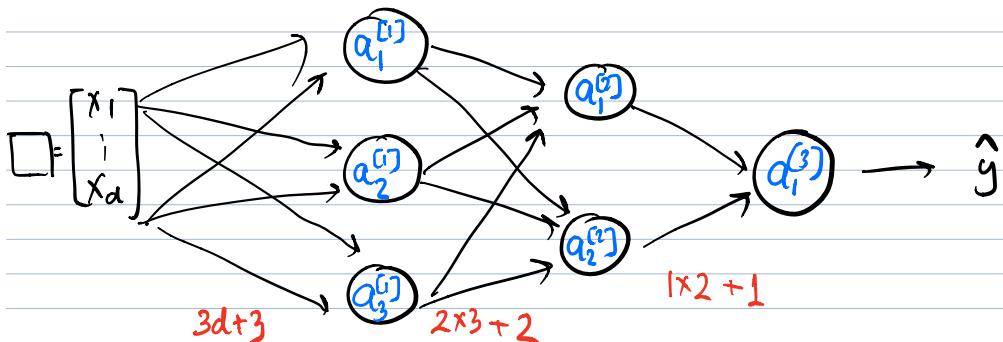
Rectified Linear Unit

Modified
Loss fn

$$\|\hat{y} - y\|_1$$

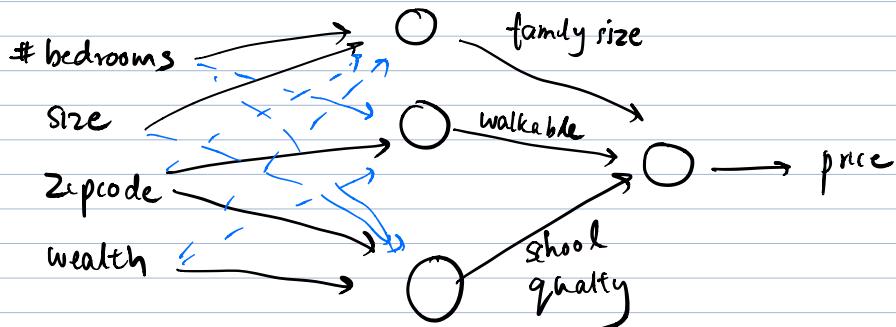
$$\|\hat{y} - y\|_2^2$$

② Neural Networks



input layer, output layer
hidden layer

House price prediction



Propagation equations

$$3 \times 1 - z^{(1)} = \underbrace{w^{(1)}x + b^{(1)}}_{3 \times d \quad d \times 1} \quad 3 \times 1$$

$$\underbrace{z^{(1)}}_{3 \times n} = \underbrace{w^{(1)} X + b^{(1)}}_{3 \times d \quad d \times n}$$

$$3 \times 1 - a^{(1)} = \sigma(z^{(1)})$$

Broadcasting

$$\underbrace{b^{(1)}}_{n \text{ copies}} = \left[\underbrace{b^{(1)}}_{1 \times d}, \dots, \underbrace{b^{(1)}}_{1 \times d} \right]$$

$$2 \times 1 - z^{(2)} = \underbrace{w^{(2)} a^{(1)} + b^{(2)}}_{2 \times 3 \quad 3 \times 1 \quad 2 \times 1} \quad 2 \times 1$$

$$2 \times 1 - a^{(2)} = \sigma(z^{(2)})$$

$$1 \times 1 - z^{(3)} = \underbrace{w^{(3)} a^{(2)} + b^{(3)}}_{1 \times 2 \quad 2 \times 1 \quad 1 \times 1} \quad 1 \times 1$$

$$1 \times 1 - a^{(3)} = \sigma(z^{(3)})$$

$$X = \begin{bmatrix} x^{(1)} & \dots & x^{(n)} \end{bmatrix}$$

[] - layer

() - id of example

capital - batch of examples

Optimize $w^{(1)}$ $w^{(2)}$ $w^{(3)}$ $b^{(0)}$ $b^{(2)}$ $b^{(3)}$

Define loss f^n