# Deep Learning

## Supervised Learning with non-linear models

before: $h_\theta(x) = \theta^T \phi(x)$      non-linear in $x$
Kernel                                           linear in $\theta$
Methods

Other non-linear models
eg. $h_\theta(x) = \sqrt{\theta_1^3 x + \theta_3 x_4 + \sqrt{\theta_5 x_8}}$

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

$$h_\theta(x) : \mathbb{R}^d \to \mathbb{R}$$

Cost / Loss $f^n$
$$J^{(i)}(\theta) = (y^{(i)} - h_\theta(x^{(i)}))^2 \quad \text{mean-squared loss}$$

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 $\{1 \cdots n\}$ uniformly
$$\theta := \theta - \alpha \nabla J^{(j)}(\theta)$$
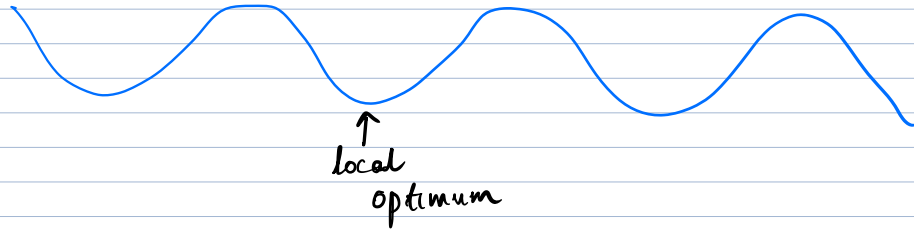
Mini-batch SGD
   — Computing $B$ gradients $\nabla^{(j_1)}(\theta) \cdots \nabla J^{(j_B)}(\theta)$
      together is faster than individual computation

Mini-batch SGD

for $i: L$ to $N_{iter}$

    Sample $B$ examples $\{j_1 \cdots j_B\}$ from
$\{1 \cdots n\}$ without replacement

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



        ↑
      local
        optimum

remaining questions

① how to define $h_\theta(x)$

② how to compute gradient

— Logistic Regression

— Neural Networks

    — computational power
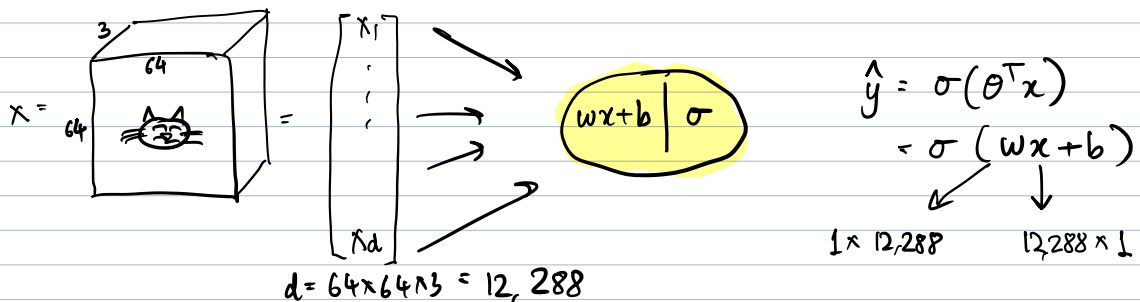
    — data available

    — algorithms

① Logistic Regression

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



$x = $    $\begin{bmatrix} x_1 \\ \vdots \\ \vdots \\ x_d \end{bmatrix}$    $\boxed{wx+b \mid \sigma}$    $\hat{y} = \sigma(\theta^\top x)$

$$= \sigma(wx+b)$$

$d = 64 \times 64 \times 3 = 12,288$      $1 \times 12,288$     $12,288 \times 1$

i) Initialize $w, b$

$\quad\quad\quad\quad\quad\quad\quad \downarrow \quad\quad \downarrow$
$\quad\quad\quad\quad\quad\quad$ weights $\quad$ bias

ii) Find optimal $w, b$

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

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

$$w := w - \alpha \frac{\partial \mathcal{L}}{\partial w}$$

$$b := b - \alpha \frac{\partial \mathcal{L}}{\partial b}$$
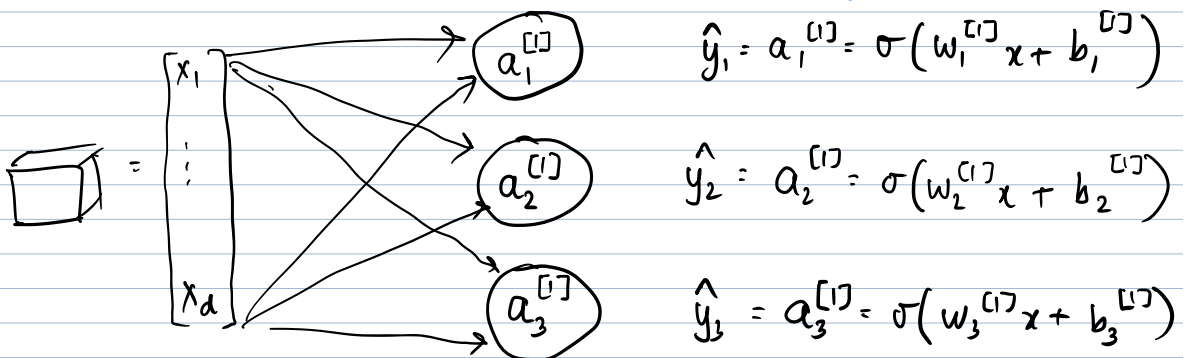
$\sigma' = \sigma(1-\sigma)$

# parameters $= 12,288 + L$

neuron $=$ linear $+$ activation

Model $=$ architecture $+$ parameters

$\quad\quad\quad\quad\quad\quad$ 1 neuron $\quad\quad\quad\quad\quad$ $w, b$

Goal 2.0 : Find cat / lion / iguana in images



$$\hat{y}_1 = a_1^{[1]} = \sigma(w_1^{[1]}x + b_1^{[1]})$$

$$\hat{y}_2 = a_2^{[1]} = \sigma(w_2^{[1]}x + b_2^{[1]})$$

$$\hat{y}_3 = a_3^{[1]} = \sigma(w_3^{[1]}x + b_3^{[1]})$$

Square brackets $[1] \equiv$ layer

Subscripts $\equiv$ identify neuron within layer

$a_2^{[1]} \leftarrow$ 1st layer

$\quad\quad \leftarrow$ 2nd neuron in 1st layer

# param $= 3(d+1)$

Dataset: images + labels

$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$ cat lion iguana
$\quad$ $\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$ lion
$\quad$ $\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$ iguana

one hot encoding

$$\begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \begin{matrix} \text{cat} \\ \text{lion} \end{matrix}$$

## Goal 3.0: add constraint: unique animal in image

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

$$\cancel{a_1^{[1]}} \quad z_1^{[1]} \qquad e^{z_1^{[1]}} / \sum_{k=1}^{3} e^{z_k^{[1]}}$$

$$\cancel{a_2^{[1]}} \quad z_2^{[1]} \qquad e^{z_2^{[1]}} / \sum_{k=1}^{3} e^{z_k^{[1]}}$$

$$\cancel{a_3^{[1]}} \quad z_3^{[1]} \qquad e^{z_3^{[1]}} / \sum_{k=1}^{3} e^{z_k^{[1]}}$$

$$\hat{y}_1 = \sigma(\underline{w_1^{[1]} x + b_1^{[1]}})$$

$$z_1^{[1]}$$

Softmax

$$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \qquad \# \text{parameters} = 3(d+1)$$

## How to train?

Loss function

$$\mathcal{L}_{3N} = -\sum_{k=1}^{3} \left[ y_k \log \hat{y}_k + (1-y_k) \log(1-\hat{y}_k) \right]$$

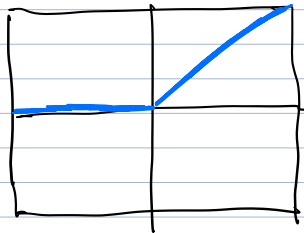$$\mathcal{L}_{CE} = -\sum_{k=1}^{3} y_k \log \hat{y}_k$$

↑ cross entropy

## $Q^n$: Instead of predicting 1/0, predict age of cat?

Options:

① Bucket ages, several neurons to predict

② Change activation $f^n$

linear f$^n$:  $f(x) = x$     linear regression

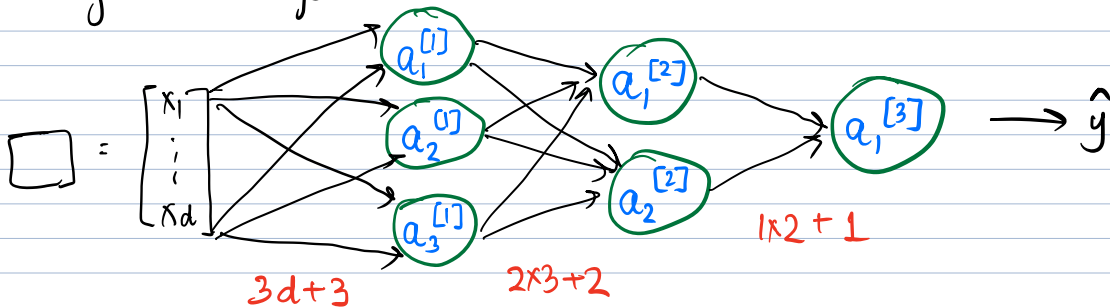$$f(x) = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases}$$

ReLU
Rectified Linear Unit

Modified    $\|\hat{y} - y\|_1$    $\|\hat{y} - y\|_2^2$
loss f$^n$
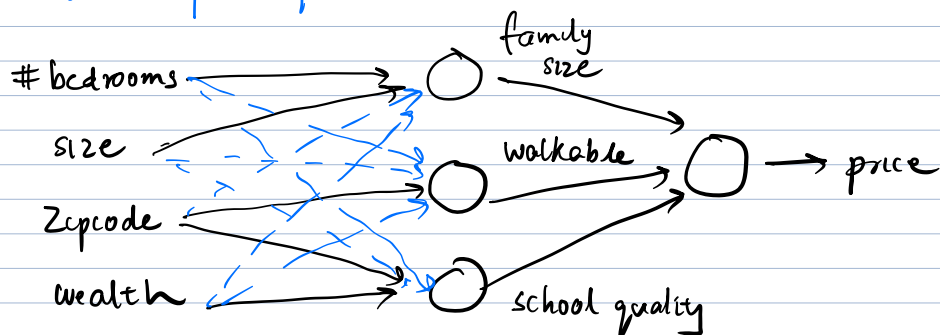
## ② Neural Networks

goal:  image $\Rightarrow$  cat  vs.  no cat



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

$a_1^{[1]}$  $a_2^{[1]}$  $a_3^{[1]}$   $a_1^{[2]}$  $a_2^{[2]}$   $a_1^{[3]}$ $\rightarrow \hat{y}$

3d+3      2×3+2      1×2+1

input layer      output layer
hidden layer

House price prediction



#bedrooms ———→  family size
size ———→  walkable ———→ price
zipcode ———→
wealth ———→  school quality

## Propagation equation

$$3\times1 - z^{[1]} = \underset{3\times d}{w^{[1]}} \underset{d\times1}{x} + \underset{3\times1}{b^{[1]}} \qquad \underset{3\times n}{Z^{[1]}} = \underset{3\times d}{w^{[1]}} \underset{d\times n}{X} + \underset{3\times1}{b^{[1]}}$$

$$3\times1 - a^{[1]} = \sigma(z^{[1]}) \qquad\qquad\qquad \text{Broadcasting}$$

$$2\times1 - z^{[2]} = \underset{2\times3}{w^{[2]}} \underset{3\times1}{a^{[1]}} + \underset{2\times1}{b^{[2]}} \qquad \tilde{b}^{[1]} = \begin{bmatrix} b^{[1]} & \cdots & b^{[1]} \end{bmatrix}$$

$$2\times1 - a^{[2]} = \sigma(z^{[2]}) \qquad \underset{\text{added to fix}}{\leftarrow} \underset{\text{omission in lecture}}{} \qquad \underbrace{\phantom{xxxxx}}_{n \text{ copies}}$$

$$1\times1 - z^{[3]} = \underset{1\times2}{w^{[3]}} \underset{2\times1}{a^{[2]}} + \underset{1\times1}{b^{[3]}}$$

$$a^{[3]} = \sigma(z^{[3]})$$

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

[ ] — layer

( ) — id of example

capital : batch of examples.

Optimize $w^{[1]}$ $w^{[2]}$ $w^{[3]}$ $b^{[1]}$ $b^{[2]}$ $b^{[3]}$

Define Loss / Cost $f^n$

## Backward Propagation