

Deep Learning

For supervised learning on non-linear models

Terminology

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 fn:

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

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

Optimization objective

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

\Rightarrow use gradient descent

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

Stochastic Gradient Descent

for $i = 1 \dots \text{niter}$

sample j from $\{1 \dots n\}$ uniformly

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

Minibatch SGD

\rightarrow Computing B gradients $\nabla J^{(j_1)}(\theta) \dots \nabla J^{(j_B)}(\theta)$
can be faster than individual computation

for $i = 1 \dots \text{niter}$

sample B examples $\{j_1 \dots j_B\}$ from $\{1 \dots n\}$ uniformly
w/o replacement

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

Note: SGD only finds local optimum

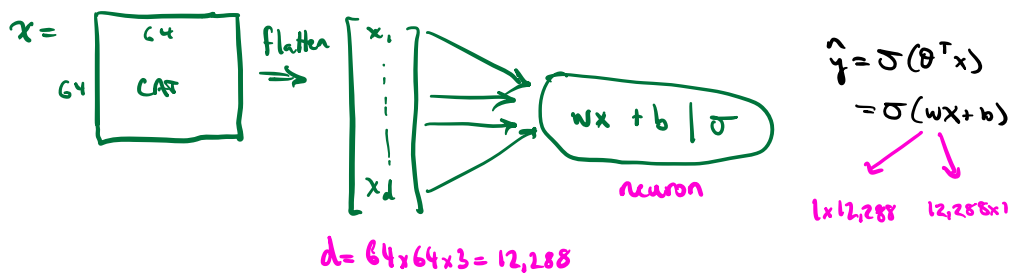
BUT: research has shown that the quality of these optima are very good compared to global

Building up a neural network

Starting from logistic regression

e.g. goal: find cats in images

1 → presence of cat
0 → absence of class



i) Initialize w, b

↓ weights ↓ bias

ii) Find optimal w, b

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

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

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

$$\sigma' = \sigma(1-\sigma)$$

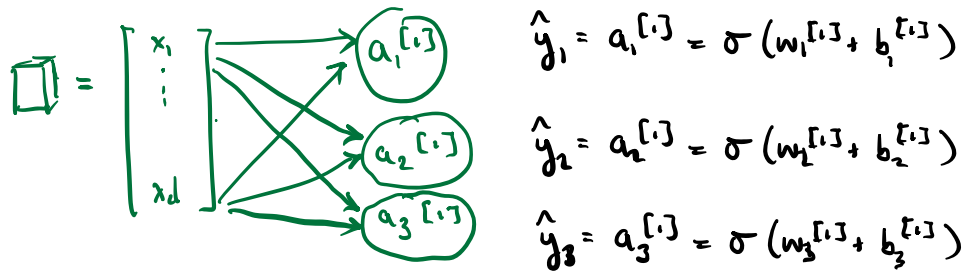
$$b' := b - \alpha \frac{\partial \mathcal{L}}{\partial b}$$

$$\# \text{parameters} = |w| + |b| = 12,288 + 1$$

define neuron linear + activation

model architecture (where neurons are, how they are arranged)
+ parameters

Goal 2.0: Find cat/lion/iguana in images



Notation: $[i] \equiv \text{layer}$

subscript: identify neuron within layer

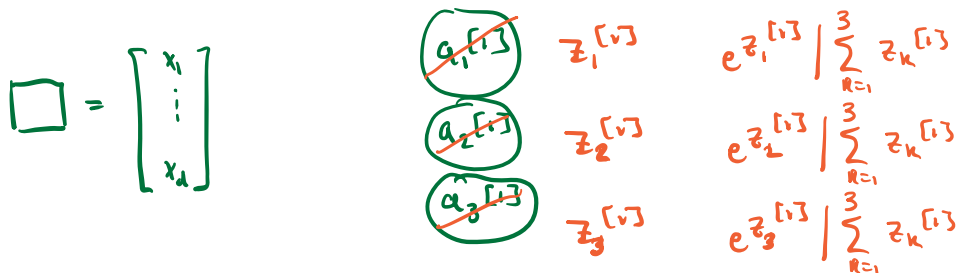
$a_2^{[1]} \leftarrow$ 1st layer
 \leftarrow 2nd neuron in layer

#params: $3(d+1)$

Dataset: images + labels

$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$ cat
 lion
 iguana one-hot encoding

Goal 3.0: add constraint: unique animal in image



$$\hat{y}_1 = \sigma(\underbrace{w_1^{[i]}x + b_1^{[i]}}_{z_1^{[i]}})$$

#params: $3(d+1)$

Training

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

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

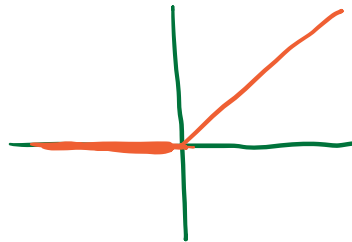
↑
cross-entropy

Predicting age of the cat

Options!

- ① Bucket ages, several neurons to predict
- ② Change activation fn

linear fn: $f(x) = x$ linear regression



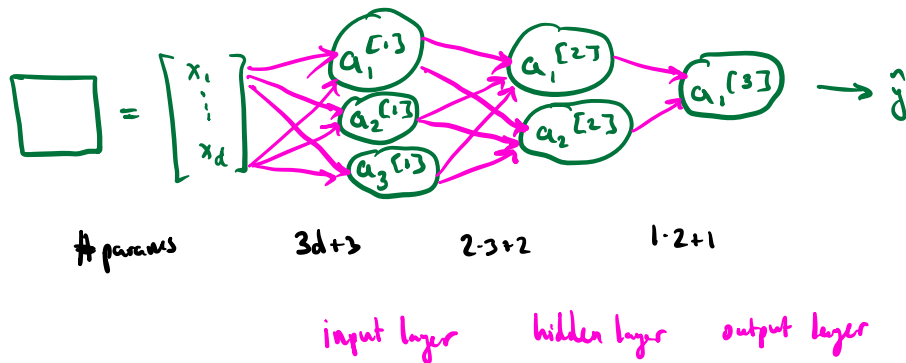
Rectified Linear Unit ReLU

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

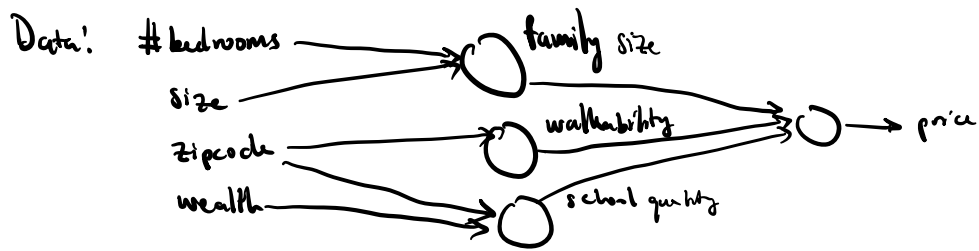
Modified loss fn $\|\hat{y} - y\|_1, \|\hat{y} - y\|_2^2$

Neural Networks, more generally

goal: image \Rightarrow cat or no cat



House price prediction



Propagation equations

$$z^{[1]} = w^{[1]} x + b^{[1]}$$

$3 \times d$ $d \times 1$ 3×1

$$a^{[1]} = \sigma(z^{[1]})$$

3x1

$$z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

2x1

$$a^{[2]} = \sigma(z^{[2]})$$

2x1

$$z^{[3]} = w^{[3]} a^{[2]} + b^{[3]}$$

1x1

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

1x1

$$z^{[1]} = w^{[1]} x + b^{[1]}$$

$3 \times n$ $3 \times d$ $d \times n$ 3×1

Broadcasting: $\tilde{b}^{[1]} = [b^{[1]} \dots b^{[1]}]$

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

[3] layer
 () id of example
 capital: batch of examples

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