

Introduction to Deep Learning

Chris Arnold, Cardiff University, July 2021

Logistic Regression

$$\hat{y} = P(y=1|x)$$

Parameters: $w \in \mathbb{R}^{n \times 1}$, $b \in \mathbb{R}$

Function $\hat{y} = \sigma(w^T x + b)$
with $\sigma(z) = \frac{1}{1 + e^{-z}}$

Loss function

goal: how far are we off?

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

if $y = 1$:

$$\mathcal{L}(\hat{y}, y) = -\log \hat{y}$$

$\left\{ \begin{array}{l} \bullet \text{ as large as possible} \\ \bullet \text{ but not } > 1 \end{array} \right.$

if $y = 0$:

$$\mathcal{L}(\hat{y}, y) = -\log(1 - \hat{y})$$

• as small as possible

• but not < 0

goal: minimise

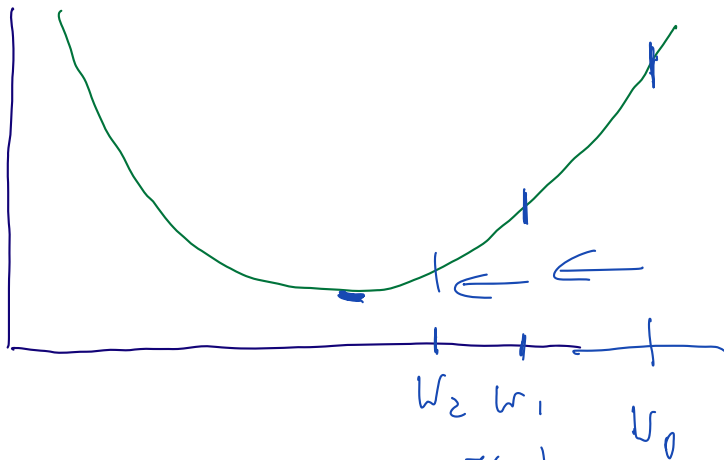
cost function

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

$$= -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}) \right]$$

Gradient Descent Intuition

Goal: $\min J(w, b)$



Repeat:

$$w := w - \alpha \frac{dJ(w)}{dw}$$

\swarrow
learning rate $\underbrace{\hspace{1cm}}$
 der. of J
 wrt w

$$w_1 = w_0 - \alpha \frac{dJ(w_0)}{dw}$$

$$w_2 = w_1 - \alpha \frac{dJ(w_1)}{dw}$$

\vdots

General Case:

$J(w, b)$:

$$w := w - \alpha \frac{dJ(w, b)}{dw}$$

$$b := b - \alpha \frac{dJ(w, b)}{db}$$

Logistic Regression Backpropagation

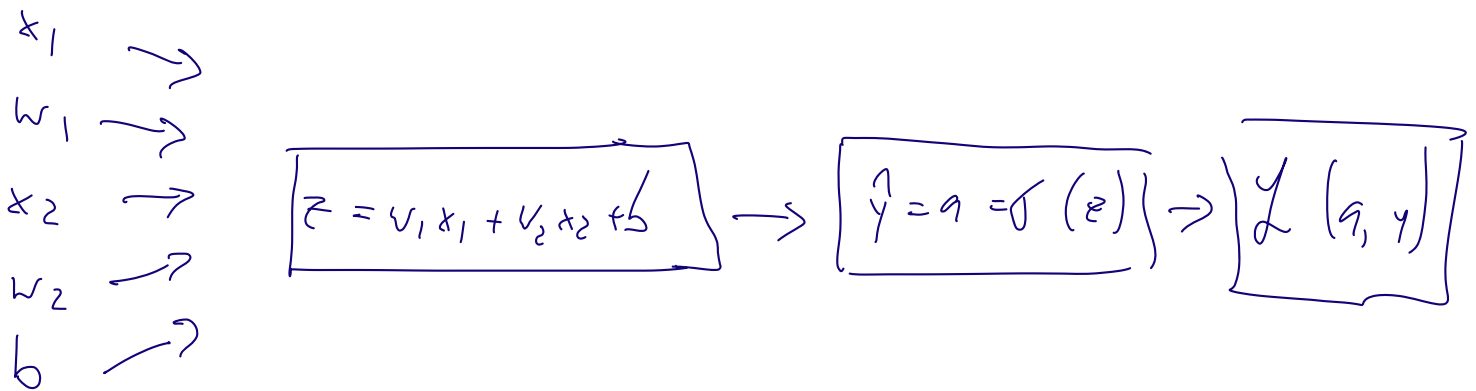
Goal: update values

Refresh

$$z = w^T x + b$$

$$\hat{y} = a = \sigma(z) = \frac{1}{1 + e^{-z}}$$

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



$$\frac{d}{da} \mathcal{L}(a, y) = -\frac{y}{a} - \frac{1-y}{1-a} \quad (1)$$

$$\frac{da}{dz} = a(1-a) \quad (2)$$

$$\frac{d\mathcal{L}}{da} \cdot \frac{da}{dz} = \left(-\frac{y}{a} - \frac{1-y}{1-a}\right) \cdot a(1-a) = a - y \quad (3)$$

$$^{\wedge}dw_1^{\wedge}: \frac{d\mathcal{L}}{dw_1} = x_1 \cdot dz$$

$$^{\wedge}dw_2^{\wedge}: \frac{d\mathcal{L}}{dw_2} = x_2 \cdot dz$$

$$^{\wedge}db^{\wedge}: \frac{d\mathcal{L}}{db} = dz$$

Updating

$$w_1 = w_1 - \alpha dw_1$$

$$w_2 = w_2 - \alpha dw_2$$

$$b = b - \alpha db$$

Neural Networks

Notation

activation a

$$a^{[0]} = x$$

$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{bmatrix}$$

$$a^{[2]} = \hat{y}$$

layer

weights w

$$w^{[1]} = \begin{bmatrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{bmatrix}$$

$$v^2 = \begin{bmatrix} \text{---} \end{bmatrix}$$

bias b

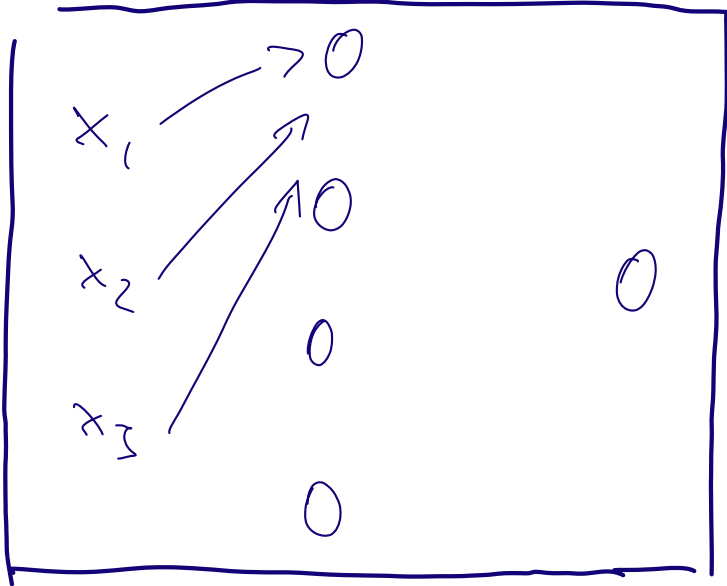
$$b^{[1]} = \begin{bmatrix} - \\ - \\ - \\ - \end{bmatrix}$$

$$b^{[2]} = [-]$$

Shallow Neural Network:

I Forward Propagation

— Zoom —



Layer 1

$$z_1^{[1]} = w_1^{[0]T} x + b_1^{[1]}$$

$$a_1^{[1]} = \sigma(z_1^{[1]})$$

$$z_2^{[1]} \dots$$

$$a_2^{[1]} \dots$$

\vdots

Vectorising

goal: compute z as a vector

$$\begin{bmatrix} w_1^{[0]T} \\ w_2^{[0]T} \\ w_3^{[0]T} \\ w_4^{[0]T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

$W^{[1]}$ $z^{[1]}$

$$\Rightarrow a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ \vdots \\ a_4^{[1]} \end{bmatrix} = \sigma(z^{[1]})$$

calculating

$$z^{[1]} = w^{[1]} a^{[0]} + b^{[1]}$$

$(4,1) \quad (4,3) \quad (3,1) \quad (4,1)$

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

$(4,1) \quad (4,1)$

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

$(1,1) \quad (1,4) \quad (4,1) \quad (1,1)$

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

$(1,1) \quad (1,1)$

II How wrong is it?

Cost function:

$$J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(\hat{y}, y)$$

III Backprop the error (skip in class)

Forward Prop:

$$z^{[1]} = w^{[1]} a^{[0]} + b^{[1]}$$

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

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

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

g can be any
activation function

Backprop

$$d z^{[2]} = a^{[2]} - y$$

$$d v^{[2]} = d z^{[2]} a^{[1]T} \quad *_2$$

$$d b^{[2]} = d z^{[2]} \quad *_1$$

$$d z^{[1]} = W^{[2]T} d v^{[2]} \odot g^{[1]'}(z^{[1]})$$

element wise product

$$d w^{[1]} = d z^{[1]} x^T \quad *_2$$

$$d b^{[1]} = d z^{[1]} \quad *_1$$

compare

compare

IV Update

$$w^{[1]} = w^{[1]} - \alpha d w^{[1]}$$

$$b^{[1]} = b^{[1]} - \alpha d b^{[1]}$$

$$w^{[2]} = w^{[2]} - \alpha d w^{[2]}$$

$$b^{[2]} = b^{[2]} - \alpha d b^{[2]}$$

Activation functions

Sigmoid

$$a = \frac{1}{1 + e^{-z}}$$

tanh

$$a = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

relu:

$$a = \max(0, z)$$

leaky relu:

$$a = \max(\alpha z, z)$$

Why activation functions?

$$a^{[0]} = z^{[1]} = W^{[1]} x + b^{[1]}$$

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

$$a^{[2]} = \underbrace{W^{[2]} W^{[1]}}_{W'} x + \underbrace{W^{[2]} b^{[1]} + b^{[2]}}_{b'}$$

L2 regularisation

log reg

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|_2^2$$

$$\|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w$$

n_x is number of cases

NN

$$J(w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_2^2$$

$$\|w^{[l]}\|_2^2 = \sum_{i=1}^{n^{[l-1]}} \sum_{j=1}^{n^{[l]}} (w_{ij}^{[l]})^2$$

Backprop

$$d w^{[l]} = \text{from backprop} + \frac{\lambda}{m} w^{[l]}$$

Updating

$$w^{[l]} := w^{[l]} - \alpha d w^{[l]}$$

$$w^{[l]} := w^{[l]} - \alpha \left[\text{from bp} + \frac{\lambda}{m} w^{[l]} \right]$$

$$w^{[l]} := w^{[l]} - \alpha (\text{from bp}) - \underbrace{\frac{\alpha \lambda}{m} w^{[l]}}$$

extra penalty
depends on

- α

- λ

- m

- but also on $w^{[l]}$