# Lecture 23: Gradient Descent, Forward and Backward Propagation

BIOS635

04/14/2020

**Class notes** for online teaching due to COVID-19.
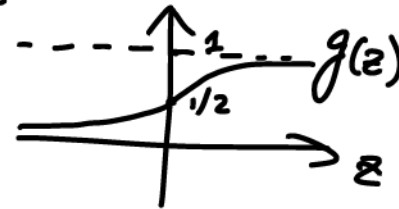
Given $x \in \mathbb{R}^p$, want $\hat{y} = P(y=1|x)$   $0 \le \hat{y} \le 1$
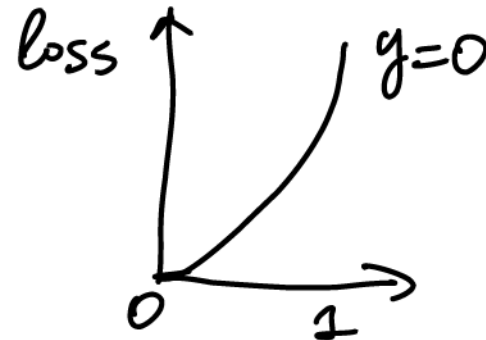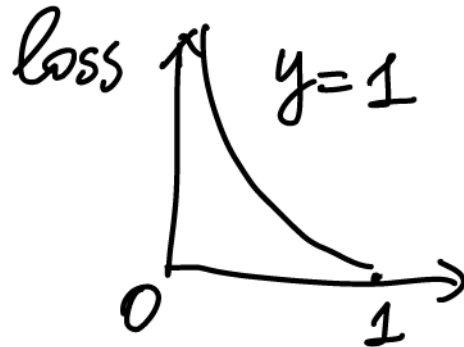
parameters $w \in \mathbb{R}^p$, $b \in \mathbb{R}$

Output   $\hat{y} = g(\underbrace{w^T x + b}_{z})$

$$g(z) = \frac{1}{1+e^{-z}}$$

Loss function: $\mathcal{L}(\hat{y}, y) = \begin{cases} -\log(\hat{y}) & \text{if } y=1 \\ -\log(1-\hat{y}) & \text{if } y=0 \end{cases}$

loss    $y=1$        loss    $y=0$

$$\mathcal{L}(\hat{y}, y) = -\underset{\text{Actual}}{y} \log \underset{\text{predicted}}{\hat{y}} - \underset{\text{Actual}}{(1-y)} \log \underset{\text{Predicted}}{(1-\hat{y})}$$ "binary Loss"

$L_2$ norm is non-convex and is hard to optimise

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

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

Want to find $w$ and $b$ that minimizes $J(W,b)$.
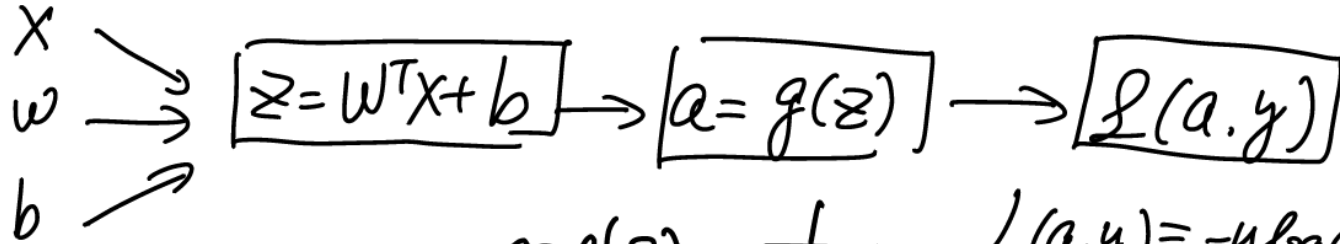Use gradient descent!

$$W := W - \alpha \, dw = W - \alpha \boxed{\frac{\partial J(W,b)}{\partial W}} \Rightarrow dw$$

$$b := b - \alpha \, db = b - \alpha \boxed{\frac{\partial J(W,b)}{\partial b}} \Rightarrow db$$

"learning rate"

## Logistic regression gradient descent

$$X$$
$$w \longrightarrow \boxed{z = W^T X + b} \longrightarrow \boxed{a = g(z)} \longrightarrow \boxed{\ell(a, y)}$$
$$b$$

$$a = g(z) = \frac{1}{1 + e^{-z}} \qquad L(a, y) = -y \log a - (1-y) \log (1-a)$$

non-linear transformation          Loss function

$$"da" = \frac{d\ell}{da} = -\frac{y}{a} + \frac{1-y}{1-a}$$

$$"dz" = \frac{dL}{dz} = \frac{dL}{da} \cdot \frac{da}{dz} = \left(-\frac{y}{a} + \frac{1-y}{1-a}\right)\left(\frac{1}{(1+e^{-z})^2} e^{-z}\right)$$

$$= \left(-\frac{y}{a} + \frac{1-y}{1-a}\right) a(1-a)$$

$$= a - y$$

$$"dw" = \frac{dL}{dw} = \frac{dL}{dz} \cdot \frac{dz}{dw} = dz \, X$$

$$"db" = \frac{dL}{db} = \frac{dL}{dz} \cdot \frac{dz}{db} = dz.$$

On m samples,

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

$$dz = a - y$$
$$dw = \frac{1}{m} X dz^T$$
$$db = \frac{1}{m} np.sum(dz)$$

$y \in \mathbb{R}^{1 \times m}, a \in \mathbb{R}^{1 \times m}, z \in \mathbb{R}^{1 \times m}$
$X \in \mathbb{R}^{P \times m}, W \in \mathbb{R}^{P}, b \in \mathbb{R}$

Putting together,

FORWARD propogation
$$\begin{bmatrix} Z = W^T X + b \\ A = g(z) \end{bmatrix}$$
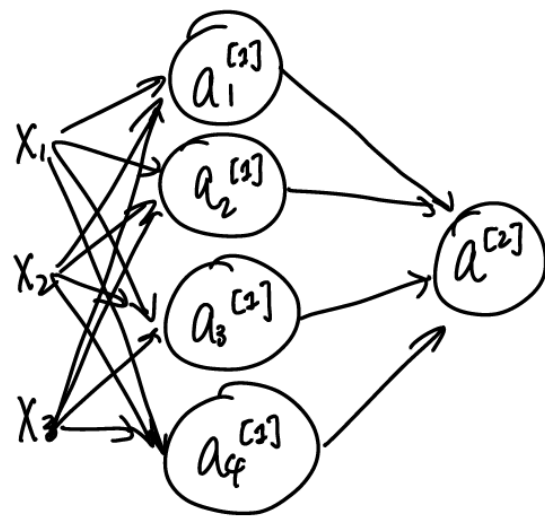
BACKWARD propagation
$$\begin{bmatrix} dz = A - Y \\ dw = \frac{1}{m} X dz^T \\ db = \frac{1}{m} np.sum(dz) \end{bmatrix}$$

$$w := w - \alpha \, dw$$
$$b := b - \alpha \, db$$

iterate till convergence

# Perceptron: one hidden-layer neural network



$$a^{[0]} = X \qquad a^{[1]} \qquad a^{[2]} = \hat{y}$$

Input layer    Hidden layer    Output layer

| FORWARD PROPAGATION |
|---|

$$\underset{n[1] \times m}{Z^{[1]}} = \underset{n[1] \times n[0]}{W^{[1]}} \underset{n[0] \times m}{A^{[0]}} + \underset{n[1] \times 1}{b^{[1]}}$$

$$\underset{n[1] \times m}{A^{[1]}} = g^{[1]}(\underset{n[1] \times m}{Z^{[1]}})$$

$$\underset{n[2] \times m}{Z^{[2]}} = \underset{n[2] \times n[1]}{W^{[2]}} \underset{n[1] \times m}{A^{[1]}} + \underset{n[2] \times 1}{b^{[2]}}$$

$$\underset{n[2] \times m}{A^{[2]}} = g^{[2]}(\underset{n[2] \times m}{Z^{[2]}})$$

First node

$$Z_1^{[1]} = W_{11}^{[1]} \cdot X_1 + W_{12}^{[1]} X_2 + W_{13}^{[1]} X_3 + b_1^{[1]}$$

F.B.t layer

Second node

$$Z_2^{[1]} = W_{21}^{[1]} X_1 + W_{22}^{[1]} X_2 + W_{23}^{[1]} X_3 + b_2^{[1]}$$

node 1

$$\vdots$$

Final node

$$Z_4^{[1]} = W_{41}^{[1]} X_1 + W_{42}^{[1]} X_2 + W_{43}^{[1]} X_3 + b_4^{[1]}$$

node 1

- parameters:

$$W^{[1]} \in \mathbb{R}^{n[1] \times n[0]}, \quad b^{[1]} \in \mathbb{R}^{n[1]}$$

$$W^{[2]} \in \mathbb{R}^{n[2] \times n[1]}, \quad b^{[2]} \in \mathbb{R}^{n[2]}$$

- Activation functions (all nonlinear)

(i) Sigmoid $g(z) = \frac{1}{1+e^{-z}}$

$g'(z) = a(1-a)$

(ii) hyperbolic tangent $g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$

$g'(z) = 1 - a^2$

(iii) Rectified linear unit (ReLU)

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

$g'(z) = \begin{cases} 0, & z < 0 \\ 1, & z \geq 0 \end{cases}$

$$a^{[2](i)}$$

$$\boxed{\text{BACKWARD PROPAGATION}}$$

Cost function: $J(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}) = \frac{1}{m}\sum_{i=1}^{m}\mathcal{L}(\hat{y}^{(i)}, y^{(i)})$
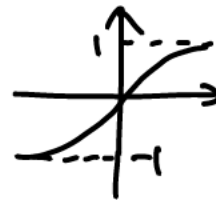
Repeat $\{$ Compute prediction $(\hat{y}^{(i)}, i=1, \cdots, m)$

$$dW^{[1]} = \frac{\partial J}{\partial W^{[1]}}, \quad db^{[1]} = \frac{\partial J}{\partial b^{[1]}}, \cdots$$

$$W^{[1]} = W^{[1]} - \alpha dW^{[1]}, \cdots$$

$$b^{[1]} = b^{[1]} - \alpha db^{[1]}, \cdots \quad \}$$

$$\underset{n^{[2]}\times m}{dZ^{[2]}} = \underset{n^{[2]}\times m}{A^{[2]}} - \underset{n^{[2]}\times m}{Y}$$

$$\underset{n^{[2]}\times n^{[1]}}{dW^{[2]}} = \frac{1}{m} \underset{n^{[2]}\times m}{dZ^{[2]}} \underset{n^{[1]}\times m}{A^{[1]T}}$$

$$\underset{n^{[2]}\times 1}{db^{[2]}} = \frac{1}{m} np.sum(\underset{n^{[2]}\times m}{dZ^{[2]}}, axis=1, keepdims=True)$$

$$\underset{n^{[1]}\times m}{dZ^{[1]}} = \underset{n^{[1]}\times m}{dA^{[1]}} * \underset{n^{[1]}\times m}{g^{[1]'}(Z^{[1]})} = \underset{n^{[1]}\times n^{[2]}}{W^{[2]T}} \underset{n^{[2]}\times m}{dZ^{[2]}} * \underset{n^{[1]}\times m}{g^{[1]'}(Z^{[1]})}$$

$$\underset{n^{[1]}\times n^{[0]}}{dW^{[1]}} = \frac{1}{m} \underset{n^{[1]}\times m}{dZ^{[1]}} \underset{m\times n^{[0]}}{X^T}$$

$$\underset{n^{[1]}\times 1}{db^{[1]}} = \frac{1}{m} np.sum(\underset{n^{[1]}\times m}{dZ^{[1]}}, axis=1, keepdims=True)$$

## Deep neural network

$$Z^{[1]} = W^{[1]} A^{[0]} + b^{[1]}$$

$$A^{[1]} = g^{[1]}(Z^{[1]})$$

$$Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$

$$A^{[2]} = g^{[2]}(Z^{[2]})$$

$$\vdots$$

$$A^{L} = g^{(L)}(Z^{[L]}) = \hat{Y}$$

$$dZ^{[L]} = A^{[L]} - Y$$

$$dW^{[L]} = \frac{1}{m} dZ^{[L]} A^{[L-1]T}$$

$$db^{[L]} = \frac{1}{m} np.sum(dZ^{[L]}, axis=1, keepdims=True)$$

$$dZ^{[L-1]} = W^{[L]T} dZ^{[L]} * g^{[L-1]'}(Z^{[L-1]})$$

$$\vdots$$

$$dZ^{[1]} = W^{[2]T} dZ^{[2]} * g^{[1]'}(Z^{[1]})$$

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} A^{[0]T}$$

$$db^{[1]} = \frac{1}{m} np.sum(dZ^{[1]}, axis=1, keepdims=True)$$

parameters : $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, W^{[3]}, b^{[3]}, \ldots$

hyperparameters : $\begin{cases} \text{learning rate } \alpha \\ \text{\# iterations} \\ \text{\# hidden layers } L \\ \text{\# hidden units } n^{[1]}, n^{[2]}, \ldots \\ \text{choice of activation function} \\ \ldots \end{cases}$

## Regularization

$$J(W^{[1]}, b^{[1]}, \cdots, W^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{\ell=1}^{L} \| W^{[\ell]} \|_F^2$$

$$W^{[\ell]} \in \mathbb{R}^{n^{[\ell]} \times n^{[\ell-1]}}, \quad \| W^{[\ell]} \|_F^2 = \sum_{i=1}^{n^{[\ell]}} \sum_{j=1}^{n^{[\ell-1]}} (W_{ij}^{[\ell]})^2$$

"Frobenius norm" (Euclidean norm of a matrix)

$$dW^{[\ell]} = (\text{from backprop}) + \frac{\lambda}{m} W^{[\ell]}$$

$$W^{[\ell]} := W^{[\ell]} - \alpha \left[ (\text{from backprop}) + \frac{\lambda}{m} W^{[\ell]} \right]$$

$$= \left(1 - \frac{\alpha\lambda}{m}\right) W^{[\ell]} - \alpha (\text{from backprop})$$

$$<1 : \text{"weight decay" for } L2 \text{ norm.}$$

Other regularization methods
- Dropout : cannot rely on any one feature
- Early stopping :



dev. set error

training set error

#iterations

$w \approx 0$     midsize $\|w\|_F^2$     large $w$