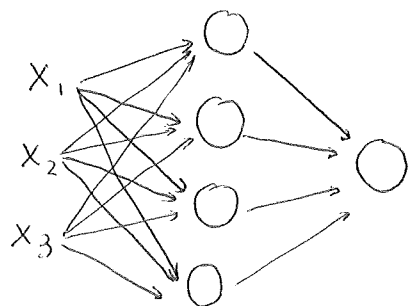


Feed forward Neural Network



Input layer Hidden Layer Output Layer

First neuron:

$$Z_1 = W_{11} x_1 + W_{12} x_2 + W_{13} x_3 + b_1$$

$$= w_1 x + b_1$$

w_1 : weight, b_1 : bias.

$$A_1 = \sigma(Z_1)$$

σ : activation function.

Hidden layer

Matrix notation:

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

$\uparrow \quad \quad \uparrow \quad \uparrow \quad \quad \uparrow$
 4x1 vector 4x3 matrix 3x1 vector 4x1 vector

superscript [1]: number of layer

$$A^{[1]} = \sigma(Z^{[1]}) \text{ component wise.}$$

\uparrow
4x1 vector

Output layer

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

$\uparrow \quad \quad \uparrow \quad \quad \uparrow \quad \quad \uparrow$
 real number 1x4 matrix 4x1 vector real number

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

As a result.

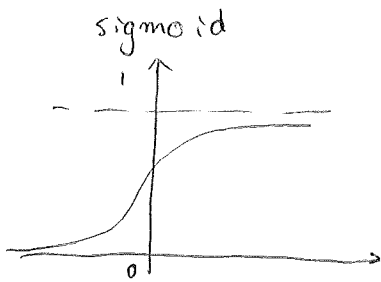
$$y = A^{[2]} = f(x_1, x_2, x_3)$$

\uparrow
composition of Linear functions and nonlinear activation functions

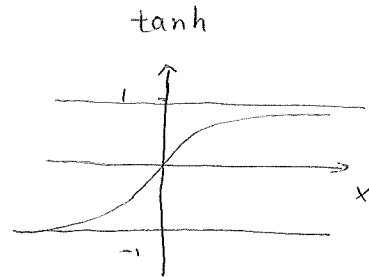
Activation function.

$$\text{sigmoid}(x) = \frac{1}{1 + e^{-x}}$$

(soft max)



$$\frac{d}{dx} \text{sigmoid}(x) = \frac{1}{1 + e^{-x}} \left(1 - \frac{1}{1 + e^{-x}} \right)$$



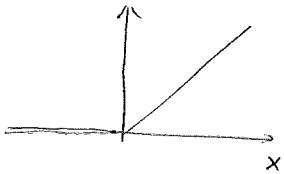
$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$\frac{d}{dx} \tanh(x) = 1 - (\tanh(x))^2$$

- $\tanh(x)$ is a scaled version of $\text{sigmoid}(x)$
- good for centering data.
- better than sigmoid
- has "vanishing gradient problem"

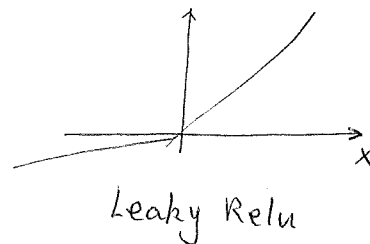
$$\text{Relu}(x) = \max(0, x)$$

Relu



- default choice.

$$\frac{d}{dx} \text{Relu}(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x < 0 \end{cases}$$

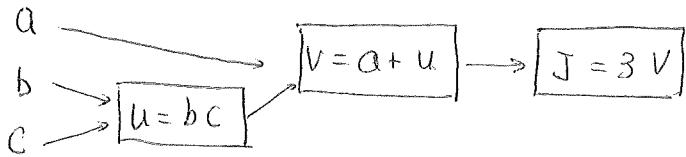


Leaky Relu

$$\text{derivative} = \begin{cases} 1 & \text{if } x > 0 \\ 0.01 & \text{if } x < 0 \end{cases}$$

Training NN

• Backward propagation.



sensitivity of J with respect a :

$$\frac{\partial J}{\partial a} = \frac{\partial J}{\partial V} \frac{\partial V}{\partial a} \quad (\text{chain rule})$$

with respect to b .

$$\frac{\partial J}{\partial b} = \frac{\partial J}{\partial V} \frac{\partial V}{\partial u} \frac{\partial u}{\partial b}$$

if J is the objective function, in gradient descend, we need

$$a_{n+1} = a_n - \underset{\substack{\uparrow \\ \text{learning rate}}}{\partial} \frac{\partial J}{\partial a}$$

• Vectorization.

$$Z = W X + b$$

$\uparrow \quad \uparrow \quad \nwarrow$
 matrix vector vector

$$Z = \text{np.dot}(w, x) + b$$

$$V = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \rightarrow u = \begin{bmatrix} e^{v_1} \\ \vdots \\ e^{v_n} \end{bmatrix}$$

$$u = \text{np.zeros}(n, 1)$$

for i in range(n):

$$u[i] = \text{math.exp}(V[i])$$

import numpy as np

$$u = \text{np.exp}(V)$$

much faster!

Logistic regression

$$X = \begin{bmatrix} \dots & x^{(1)} & \dots \\ \vdots & \vdots & \vdots \\ \dots & x^{(n)} & \dots \end{bmatrix} \quad (n, p) \text{ matrix}$$

$$\begin{bmatrix} z^{(1)} \\ \vdots \\ z^{(n)} \end{bmatrix} = W X + \begin{bmatrix} b \\ \vdots \\ b \end{bmatrix}$$

$$(1) \quad Z = \text{np.dot}(W, X) + b \quad \text{response variable.}$$

$$A = \underset{\substack{\uparrow \\ \text{sigmoid.}}}{\sigma}(Z) \quad A = \begin{bmatrix} a^{(1)} \\ \vdots \\ a^{(n)} \end{bmatrix} \quad \downarrow \quad Y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

$$(2) \quad \text{denote } dZ = A - Y$$

we have seen from note before that

$$\frac{\partial J}{\partial W} = \frac{1}{n} X^T dZ$$

// denote

dW

$$dW = 0.$$

$$dW_{+} = x^{(1)} dZ^{(1)}$$

$$dW_{+} = x^{(2)} dZ^{(2)}$$

\vdots

$$dW_{/} = n$$

$$\frac{\partial J}{\partial b} = \frac{1}{n} \sum_{i=1}^n dZ^{(i)}$$

// denote

db

$$db = 0$$

$$db_{+} = dZ^{(1)}$$

$$db_{+} = dZ^{(2)}$$

\vdots

$$db_{/} = n.$$

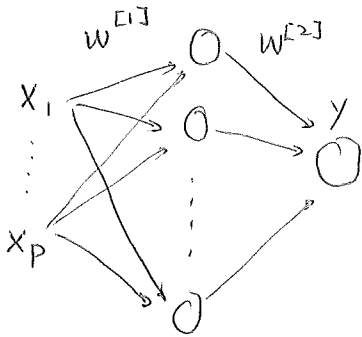
$$(3) \quad dW = \frac{1}{n} \text{np.dot}(X, T, dZ) \quad db = \frac{1}{n} \text{np.sum}(dZ)$$

$$(4) \quad w_{-} = -\partial dW \quad b_{-} = -\partial db.$$

Gradient descent for NN

5

$$J = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) \quad : \text{ Loss function}$$



sigmoid is chosen as activation function for the output layer.

Forward: $Z^{[1]} = W^{[1]} X + b^{[1]}$

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

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

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

Backward propagation:

$$dZ^{[2]} = A^{[2]} - Y$$

$$dW^{[2]} = \frac{1}{n} (A^{[1]})^T dZ^{[2]}$$

$$db^{[2]} = \frac{1}{n} np. \text{sum}(dZ^{[2]})$$

} Same as Logistic regression with $A^{[1]}$ as X .


$$dZ^{[1]} = W^{[2]} dZ^{[2]} \underset{\substack{\uparrow \\ \text{element product}}}{*} (g^{[1]})'(Z^{[1]})$$

$$\left(\begin{array}{l} \text{chain rule} \\ \frac{\partial J}{\partial Z^{[1]}} = \frac{\partial J}{\partial Z^{[2]}} \underbrace{\frac{\partial Z^{[2]}}{\partial A^{[1]}}}_{W^{[2]}} \underbrace{\frac{\partial A^{[1]}}{\partial Z^{[1]}}}_{(g^{[1]})'(Z^{[1]})} \end{array} \right)$$

$$dW^{[1]} = \frac{1}{n} X^T dZ^{[1]}$$

$$db^{[1]} = \frac{1}{n} np. \text{sum}(dZ^{[1]})$$

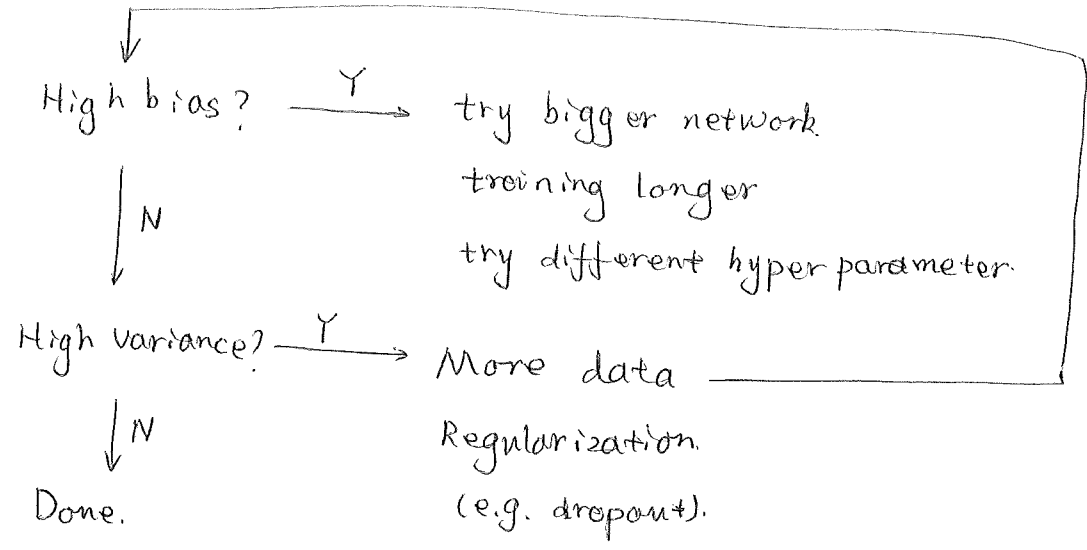
Training for NN

Data. 
 Training Validation test.

rule of thumb : 60 - 20 - 20.

Bigdata 1,000,000, - 10,000 - 10,000.

Validation is used to test different models.

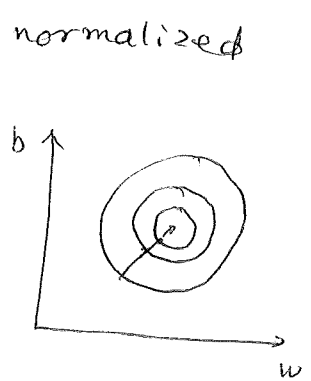
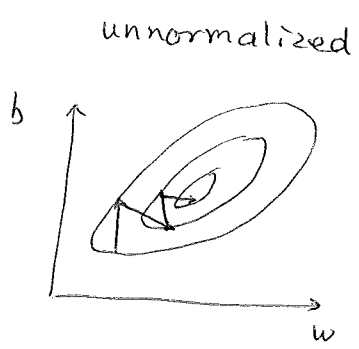


Optimization problem in NN

Setting up optimization problem:

Normalizing inputs $x \rightarrow \frac{x - \mu}{\sigma}$

why?



Exploding or vanishing gradients

deep net work; $\hat{y} = w^{[L]} w^{[L-1]} \dots w^{[1]} x$. $b^{[i]} = 0 \forall i$.

$$w^{[1]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix} \quad w^{[2]} = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$$

$$\hat{y} = \begin{bmatrix} 1.5^L & 0 \\ 0 & 1.5^L \end{bmatrix} x \quad \text{or} \quad \hat{y} = \begin{bmatrix} 0.5^L & 0 \\ 0 & 0.5^L \end{bmatrix} x$$

exploding

vanishing

random weight initialization

$$Z = w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$

~~var~~ $\{w_i\}$ i.i.d.

$$\text{Var}(w_i) = \frac{1}{n} \quad \text{larger } n \rightarrow \text{smaller } w_i$$

$$w^{[i]} = \text{np.random.randn(shape)} * \text{np.sqrt}\left(\frac{1}{n^{[i]}}\right)$$

Optimization

mini-batch gradient descent

$$x = \underbrace{[x^{(1)}, x^{(2)}, \dots, x^{(n)}]}_{\substack{X^{[1]} \\ 1000}} \underbrace{[x^{(n+1)}, \dots, x^{(2n)}]}_{X^{[2]}} \underbrace{[x^{(2n+1)}, \dots, x^{(m)}]}_{\substack{X^{[5000]} \\ 5,000,000}}$$

Calculate gradient using mini batch.

repeat { going through data multiple times.

for $t = 1, \dots, 5000$

1 mini batch. {

Forward prop on $X^{[t]}$

For i in each layer

$$Z^{[i]} = W^{[i]} X^{[t]} + b^{[i]}$$

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

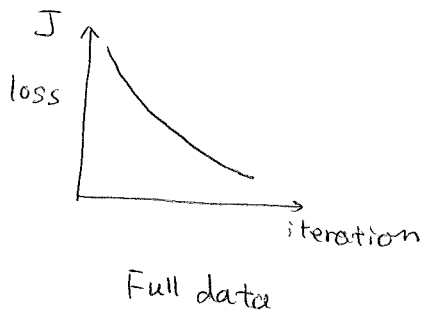
compute cost

$$J^{[t]} = \frac{1}{1000} \sum_{j=1}^{1000} \mathcal{L}(\hat{y}^{(j)}, y^{(j)}) + \frac{\lambda}{2} \frac{1}{1000} \sum \|W^{[i]}\|^2$$

Gradient descent to upgrade all $w^{[i]}$

} 1 epoch.

Batch gradient descent



Size of mini-batch.

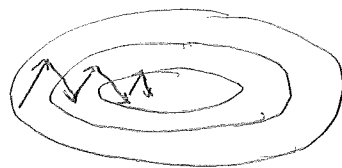
mini batch size = Full data size : Gradient descent

mini batch size = 1 : Stochastic gradient descent

In practice : between 1 and Full data size.

typically 64, 128, 256, 512

Gradient descent with momentum



On iteration t ,

compute dw, db on mini-batch.

$$V_{dw} = \beta V_{dw} + (1 - \beta) dw$$

$$V_{db} = \beta V_{db} + (1 - \beta) db \quad \beta = 0.9$$

$$w = w - \partial V_{dw} \quad b = b - \partial V_{db}$$

Oscillation in vertical direction almost cancel each other.

We want faster descent
on horizon direction and
slower descent on vertical
direction.

RMS prop

On iteration t ,

compute dw, db on mini-batch.

$$S_{dw} = \beta S_{dw} + (1 - \beta) (dw)^2$$

$$S_{db} = \beta S_{db} + (1 - \beta) (db)^2$$

$$w = w - \partial \frac{dw}{\sqrt{S_{dw} + \epsilon}}$$

$$b = b - \partial \frac{db}{\sqrt{S_{db} + \epsilon}}$$

if db is large, S_{db} is large,
the ratio dampen updating in b .

Adam (adaptive momentum estimation)

$$V_{dw} = 0, S_{dw} = 0, V_{db} = 0, S_{db} = 0.$$

On iteration t ,

compute dw, db using mini-batch.

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) dw$$

$$V_{db} = \beta_1 V_{db} + (1 - \beta_1) db$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) (dw)^2$$

$$S_{db} = \beta_2 S_{db} + (1 - \beta_2) (db)^2$$

$$V_{dw}^{\text{corrected}} = V_{dw} / (1 - \beta_1^t)$$

$$V_{db}^{\text{corrected}} = V_{db} / (1 - \beta_1^t)$$

$$S_{dw}^{\text{corrected}} = S_{dw} / (1 - \beta_2^t)$$

$$S_{db}^{\text{corrected}} = S_{db} / (1 - \beta_2^t)$$

$$w = w - \partial \frac{V_{dw}^{\text{corrected}}}{\sqrt{S_{dw}^{\text{corrected}} + \epsilon}}$$

$$b = b - \partial \frac{V_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{corrected}} + \epsilon}}$$