

CSC413 Homework 1

1.1. Verify Sort

let h_1 detect $x_1 \leq x_2 \Leftrightarrow x_2 - x_1 \geq 0$
 let h_2 detect $x_2 \leq x_3 \Leftrightarrow x_3 - x_2 \geq 0$
 let h_3 detect $x_3 \leq x_4 \Leftrightarrow x_4 - x_3 \geq 0$
 let y detect $h_1 = h_2 = h_3 = 1$

$$h = W^{(1)}x + b^{(1)} \quad W^{(1)} = \begin{bmatrix} -1 & 1 & 0 & 0 \\ 0 & -1 & 1 & 0 \\ 0 & 0 & -1 & 1 \end{bmatrix}$$

$3 \times 1 \quad 3 \times 4 \quad 4 \times 1 \quad 3 \times 1$

$$b^{(1)} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$y = W^{(2)T}h + b^{(2)} \quad W^{(2)} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

$1 \times 1 \quad 1 \times 3 \quad 3 \times 1 \quad 1 \times 1$

$$b^{(2)} = -2.5$$

1.3 Universal Approximation Theorem

1.3.1 bump function.

if $a \leq x \leq b$, return 1

otherwise return 0

* note that the final layer does not have an activation function

$$h = W_0x + b_0$$

$2 \times 1 \quad 2 \times 1 \quad 1 \times 1 \quad 2 \times 1$

$$W_0 = \begin{bmatrix} 1 & -1 \end{bmatrix}$$

$$b_0 = \begin{bmatrix} -a \\ +b \end{bmatrix}$$

$$y = W_1h + b_1$$

$1 \times 1 \quad 1 \times 2 \quad 2 \times 1 \quad 1 \times 1$

$$W_1 = \begin{bmatrix} h & 1 \end{bmatrix}^T$$

$$b_1 = -h$$

$$W_1W_0^T + W_1b_0 - b_1 = h$$

$$W_1 + W_2 - b_1 = h$$

1.2 Perfrom Sort

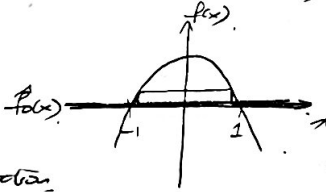
Suppose we have now 1. the verify-sort that checks for the sortedness

and

2. a permutation neural network, which takes in $[x_1, x_2, \dots, x_p]$ and output all $p!$ permutations.

then, for p in permutations:
 if verify-sort(p) == 1:
 return p .

$$1.3.2 \quad f(x) = -x^2 + 1, \quad f_0(x) = 0$$

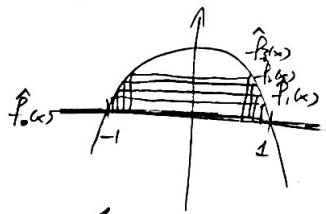


$$\text{find } \hat{f}_1(x) = f_0(x) + g(h_1, a_1, b_1, \tau)$$

$$= f_0(x) + h_1 \cdot I(a_1 \leq x \leq b_1)$$

let a_1 be a value slightly greater than -1, say $a_1 = -0.999$
 let b_1 be a value slightly smaller than 1, say $b_1 = 0.999$
 let h_1 be $f(a_1)$.

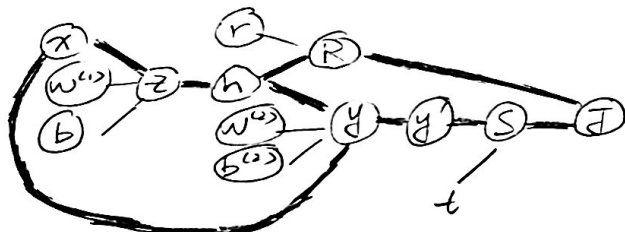
1.3.3



let a_2 be a value slightly greater than a_1 .
 let b_2 be a value slightly smaller than b_1 .
 let h_2 be $f(a_2)$
 keep doing this until a_n and b_n meets at 0

2.1 Computational Graph.

2.1.2 $\bar{J} = 1$



2.1.2, assuming $M=N$

$$\rightarrow -z = w^{(1)}x + b^{(1)}$$

$$\begin{array}{ccc} \rightarrow h = \text{ReLU}(z) & & \rightarrow R = r^T h \\ k \times 1 & & 1 \times 1 \quad 1 \times k \quad k \times 1 \end{array}$$

$$\rightarrow y = w^{(2)}h + b^{(2)} + x$$

$m \times 1$ $m \times k$ $k \times 1$ $m \times 1$ $N \times 1$

$$\rightarrow y' = \text{softmax}(y) \rightarrow S = I(t=k) y'_k$$

$$f(x) = \sqrt{V^T x}$$

2.2.1

$$v = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

$$VV^T = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 \end{bmatrix} = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 6 \\ 3 & 6 & 9 \end{bmatrix}$$

$$v^T x = \begin{bmatrix} 1 & 2 & 3 \\ 3 & 6 & 9 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} x_1 + 2x_2 + 3x_3 \\ 3x_1 + 6x_2 + 9x_3 \end{bmatrix} = f$$

$$J = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \frac{\partial f_1}{\partial x_3} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \frac{\partial f_2}{\partial x_3} \\ \frac{\partial f_3}{\partial x_1} & \frac{\partial f_3}{\partial x_2} & \frac{\partial f_3}{\partial x_3} \end{bmatrix} = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 6 \\ 3 & 6 & 9 \end{bmatrix}$$

2.2.2

time cost = n^2
memory cost = n

2.2.3

$$\underline{J} = \frac{\partial f}{\partial \underline{x}} = \underline{v} \underline{v}^T$$

$$\bar{x}^T y = (v v^T)^T y = v v^T y$$

V_T requires $\sim n$ multiply-additions

$v(vTy)$ requires n multiplies.

So time cost is linear, memory cost is also linear.

method: $J^T y = \sqrt{V}^T y$.

$$\therefore \delta \bar{y} = \sqrt{V} \bar{y} = \begin{bmatrix} 1 \\ 3 \end{bmatrix} [1 \ 2 \ 3] \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}, b = \begin{bmatrix} 6 \\ 12 \\ 18 \end{bmatrix} \quad \checkmark$$

$$\bar{U} = \frac{\partial \bar{U}^T}{\partial \bar{S}} \bar{U} = 1 \bar{U}$$

$$\bar{y}' = \frac{\partial S^T}{\partial y' S} = 1I(t > k)$$

$$\sqrt{y} = \frac{\partial y'}{\partial y} y' = \text{softmax}'(y, y')$$

$$\bar{R} = \frac{\partial \bar{I}}{\partial \bar{p}} = 1 \bar{p}$$

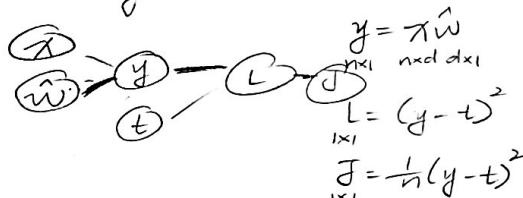
$$\bar{h} = \frac{\partial R}{\partial n} \bar{R} + \frac{\partial y}{\partial n} \bar{y} = r \bar{R} + w(2) \bar{y}$$

$$\sqrt{z} = \frac{\partial h_T}{\partial z} \frac{1}{h} = \frac{1}{h} \max(0, z) \frac{1}{z}$$

$$\checkmark \bar{x} = \frac{\partial^2 T}{\partial x^T} \bar{z} = \omega^{(1)T} \bar{z} + \bar{y}$$

3. Linear Regression.

3.1 Deriving the Gradient.



$$\bar{f} = 1.$$

$$\bar{L} = \frac{1}{n}$$

$$\bar{y} = \frac{\partial L^T}{\partial y} L = \frac{1}{n} \sum (y - t) = \frac{2}{n} (y - t)$$

$$\frac{\partial \omega}{\partial x_1} = \frac{\partial y^T}{\partial \omega} \frac{\partial y}{\partial x_1} = \frac{2}{n} x^T (y - t) = \frac{2}{n} x^T (x \omega - t)$$

3.2 Underparametrized Model.

3.2.1

$$\bar{\omega} = 0 = \underset{\text{dnn}}{\pi^T} (\underset{\text{nn}}{y} - t) = \underset{\text{dnn}}{\pi^T} (\underset{\text{nn}}{\pi} \underset{\text{nn}}{\omega^*} - t)$$

$$\frac{dw^k}{dx_j} = (-\pi^T \pi)^{-1} \pi^T e$$

3.2.2

$$\begin{aligned}\hat{w}^* &= (\pi^T \pi)^{-1} \pi^T \pi w^* \\ &= (\pi^T \pi)^{-1} \pi^T \pi w^* \\ &= w^*\end{aligned}$$

$$\text{So, } (w^T x - w^T x)^2 = 0$$

3.3 Overparametrized Model: 2D Example

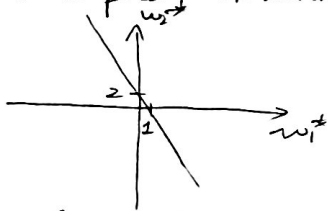
3.3.1.

$$\begin{aligned}\bar{w} = 0 &= \mathbf{x}^T(\mathbf{y} - \mathbf{t}) = \mathbf{x}^T(\mathbf{x}\hat{\mathbf{w}}^* - \mathbf{t}) \\ \mathbf{x}^T\mathbf{x}\hat{\mathbf{w}}^* &= \mathbf{x}^T\mathbf{t} \\ \mathbf{x}\mathbf{x}^T\mathbf{x}\hat{\mathbf{w}}^* &= \mathbf{x}\mathbf{x}^T\mathbf{t} \\ \mathbf{x}\hat{\mathbf{w}}^* &= \mathbf{t}\end{aligned}$$

$$\begin{bmatrix} 2 \\ 1 \end{bmatrix} \begin{bmatrix} w_1^* \\ w_2^* \end{bmatrix} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

$$2w_1^* + w_2^* = 2$$

any (w_1^*, w_2^*) satisfy this relationship is an empirical risk minimizer.



$$\text{line: } w_2^* = 2 - 2w_1^*$$

3.3.2

when $\hat{w}(0) = 0$, find direction of gradient.

$$\begin{aligned}\bar{w} &= \frac{2}{1} \begin{bmatrix} 2 \\ 1 \end{bmatrix} (0 - 2) \\ &= -4 \begin{bmatrix} 2 \\ 1 \end{bmatrix} = \begin{bmatrix} -8 \\ -4 \end{bmatrix}\end{aligned}$$

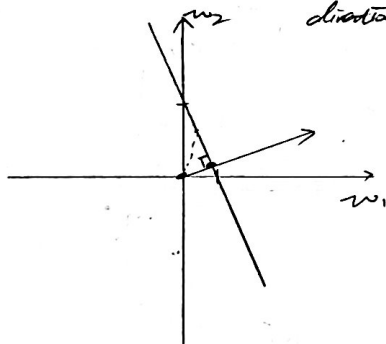
after normalization:

$$\bar{w} = \begin{bmatrix} -2/\sqrt{5} \\ -1/\sqrt{5} \end{bmatrix}, \text{ direction of gradient descent is } \begin{bmatrix} 2/\sqrt{5} \\ 1/\sqrt{5} \end{bmatrix}$$

direction doesn't change.

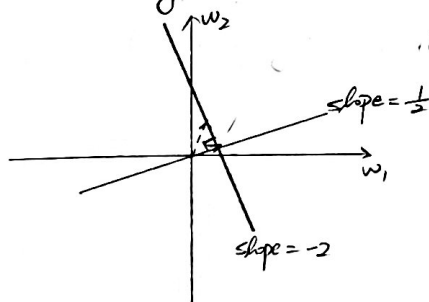
solution: find intersection of

$$\begin{aligned}w_2 &= 2 - 2w_1 \text{ and } w_2 = \frac{1}{2}w_1 \\ \frac{1}{2}w_1 &= 2 - 2w_1 \\ \frac{1}{2}w_1 + 2w_1 &= 2 \\ w_1 &= \frac{4}{5}, w_2 = \frac{2}{5}\end{aligned}$$



3.3.3.

Geometrically,



Since the slopes are negative reciprocals, the gradient descent solution is perpendicular to the line of all empirical risk minimizers.

Suppose there exists another solution on the line, using Pythagorean, $a^2 + b^2 = c^2$.

$$\begin{aligned}c^2 &= a^2 + b^2 \\ c^2 &> a^2 \\ c &> a\end{aligned}$$

3.4 Overparametrized Model: General Case.

3.4.1.

$$\bar{w} = 0 \Rightarrow \mathbf{x}\hat{\mathbf{w}}^* = \mathbf{t} \quad \text{empirical risk minimizer}$$

$$\begin{aligned}\bar{w}(0) &= \frac{2}{n} \mathbf{x}^T(\mathbf{x}\hat{\mathbf{w}}^* - \mathbf{t}) \\ &= \frac{2}{n} \mathbf{x}^T\mathbf{t} - \mathbf{t}\end{aligned}$$

gradient direction is a linear combination of rows of \mathbf{x}^T , the coefficients are the n elements in \mathbf{t}

\Rightarrow gradient is spanned by the rows of \mathbf{x}^T

find $\hat{\mathbf{w}}^*$ that ① satisfies $\mathbf{x}\hat{\mathbf{w}}^* = \mathbf{t}$

② $\hat{\mathbf{w}}^*$ spanned by rows of \mathbf{x}^T

let $\mathbf{c} \in \mathbb{R}^{n \times 1}$ be the coefficients

$$\mathbf{x}\mathbf{x}^T\mathbf{c} = \mathbf{t}$$

$$\mathbf{c} = (\mathbf{x}\mathbf{x}^T)^{-1}\mathbf{t}$$

$$\text{So, } \hat{\mathbf{w}}^* = \mathbf{x}^T(\mathbf{x}\mathbf{x}^T)^{-1}\mathbf{t}$$

3.4.2

$$\begin{aligned}(\hat{\mathbf{w}}^* - \hat{\mathbf{w}})^T \hat{\mathbf{w}}^* &= (\hat{\mathbf{w}}^{*T} - \hat{\mathbf{w}}^T) \hat{\mathbf{w}}^* \\ &= \hat{\mathbf{w}}^{*T} \hat{\mathbf{w}}^* - \hat{\mathbf{w}}^T \hat{\mathbf{w}}^* \\ &= \mathbf{t}^T (\mathbf{x}\mathbf{x}^T)^{-1} \mathbf{x} (\mathbf{x}\mathbf{x}^T)^{-1} \mathbf{t} - \hat{\mathbf{w}}^T \mathbf{x}^T (\mathbf{x}\mathbf{x}^T)^{-1} \mathbf{t} \\ &= \mathbf{t}^T (\mathbf{x}\mathbf{x}^T)^{-1} \mathbf{t} - \hat{\mathbf{w}}^T \mathbf{x}^T (\mathbf{x}\mathbf{x}^T)^{-1} \mathbf{t} \\ &= \mathbf{t}^T (\mathbf{x}\mathbf{x}^T)^{-1} \mathbf{t} - (\hat{\mathbf{w}}^T \mathbf{x}^T (\mathbf{x}\mathbf{x}^T)^{-1}) \mathbf{t} \\ &= \mathbf{t}^T (\mathbf{x}\mathbf{x}^T)^{-1} \mathbf{t} - (\mathbf{x}\mathbf{x}^T)^{-1} \mathbf{t} = 0\end{aligned}$$

So, $(\hat{\mathbf{w}}^* - \hat{\mathbf{w}})$ is orthogonal to $\hat{\mathbf{w}}^*$.

using the same reasoning of Pythagorean, the norm of $\hat{\mathbf{w}}^*$ is the smallest.



3.1 Benefit of Overparameterization,

3.1.1

code snippet:

if $d \geq n$:

$w = \mathbf{X}_{\text{expand}}^T @ \text{linAlg.inv}(\mathbf{X}_{\text{expand}} @ \mathbf{X}_{\text{expand}}) @ \mathbf{X}_{\text{expand}} @ y$

else:

$w = \text{linAlg.inv}(\mathbf{X}_{\text{expand}}^T @ \mathbf{X}_{\text{expand}}) @ \mathbf{X}_{\text{expand}}^T @ y$

return w

No, overparameterization doesn't always
lead to overfitting. ✓