# 1. Hard-Coding Networks.

## 1.1. Verify Sort

$$W^{(1)} = \begin{bmatrix} -1 & 1 & 0 & 0 \\ 0 & -1 & 1 & 0 \\ 0 & 0 & -1 & 1 \end{bmatrix} \qquad b^{(1)} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$W^{(2)} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \qquad b^{(2)} = -2.5$$

## 1.2. Perform Sort:
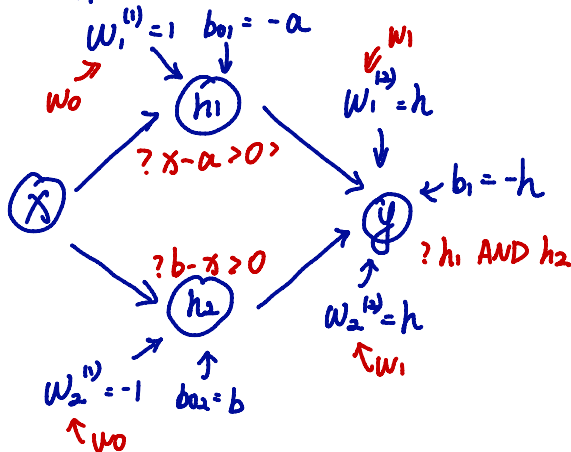
$$P \leftarrow \text{Permutation } (x_1, x_2, x_3, x_4).$$

for $p$ in $P$ :

     if verify sort $(p) > 0$:

         return $p$

## 1.3. Universal Approximation Theorem.

### 1.3.1.



$$W_1 \, a(W_0 x + b_0) + b_1 = h \, \mathbb{1} (a < x < b)$$

$$W_0 = \begin{bmatrix} 1 \\ -1 \end{bmatrix} \quad b_0 = \begin{bmatrix} -a \\ b \end{bmatrix} \quad a = \mathbb{1}(y \geq 0)$$

$$W_1 = [h, h]. \quad b_1 = -h$$

### 1.3.2.

$$\|f - \hat{f}_1\| = \int_I |f(x) - \hat{f}_1(x)| dx = \int_{-1}^1 |-x^2+1 - 0 - g(h_1, a_1, b_1, x)| dx \text{①}, \text{ suppose } [a, b] \in I$$

$$= \int_{-1}^a -x^2+1 \, dx + \int_b^1 -x^2+1 \, dx + \int_a^b |-x^2+1-h| dx$$

① ⇒

$|-x^2+1-h| > 0$

$$= \int_{-1}^1 -x^2+1 \, dx + \int_a^b -h \, dx \qquad \int_I -x^2+1 = -\frac{x^3}{3} + x \Big]_{-1}^1 = \frac{4}{3} = \|f - \hat{f}_0\|$$

$$= \frac{4}{3} - h(b-a) < \frac{4}{3}$$

$$h(b-a) > 0$$

$(-x^2+1-h) < 0$

① ⇒

$$\int_{-1}^a -x^2+1 \, dx + \int_b^1 -x^2+1 \, dx + \int_a^b x^2-1+h \, dx$$

$$= -\frac{x^3}{3} + x \Big]_{[b, 1]}^{[-1, a]} + \frac{x^3}{3} - x + hx \Big]_a^b$$

$$= -\frac{a^3}{3} + a - [\frac{1}{3}-1] + -\frac{1}{3}+1 - [\frac{b^3}{3}+b] + \frac{b^3}{3} - b + hb \quad \{\frac{a^3}{3} - a + ha\}$$

$$= -\frac{2}{3}a^3 + \frac{2}{3}b^3 + 2a - 2b + hb - ha + \frac{4}{3} < \frac{4}{3}$$
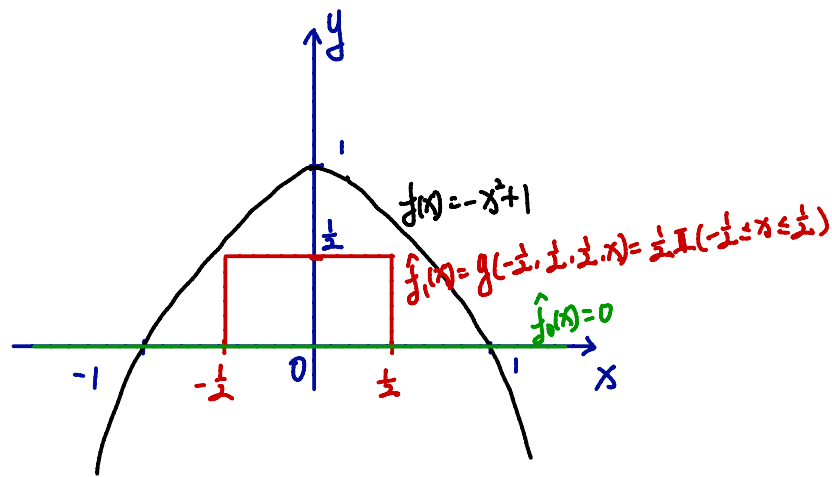
$$(b-a)(h-2) < \frac{2}{3}(a^3-b^3)$$

take $\hat{f}_1 = g(-\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, x).$

$b = \frac{1}{2}, \ a = -\frac{1}{2}, \ h = \frac{1}{2}$

$$(\frac{1}{2}+\frac{1}{2}) \cdot (\frac{1}{2}-2) < \frac{2}{3}((-\frac{1}{2})^3 - (\frac{1}{2})^3)$$

$$-\frac{3}{2} < -\frac{1}{6}$$

$\frac{1}{2} \cdot (\frac{1}{2}+\frac{1}{2}) \geq 0 \Rightarrow$ checks.

$$f(x) = -x^2 + 1$$
$$\hat{f_1}(x) = g(-\tfrac{1}{2}, \tfrac{1}{2}, \tfrac{1}{2}, x) = \tfrac{1}{2}\mathbb{I}(-\tfrac{1}{2} \le x \le \tfrac{1}{2})$$
$$\hat{f_0}(x) = 0$$

1.3.3.
$$\hat{f_1}(x) = \hat{f_0}(x) + g(-\tfrac{1}{2^N}, \tfrac{1}{2^N}, \tfrac{1}{2^N}, x).$$
$$\hat{f_2}(x) = \hat{f_1}(x) + g(-\tfrac{1}{2^{N-1}}, \tfrac{1}{2^{N-1}}, \tfrac{1}{2^{N-1}}, x)$$
$$\vdots$$
$$\hat{f_N}(x) = \hat{f_{N-1}}(x) + g(-\tfrac{1}{2}, \tfrac{1}{2}, \tfrac{1}{2}, x)$$
$$\Rightarrow$$

$$N = 3. \quad \hat{f_1}(x) = g(-\tfrac{1}{8}, \tfrac{1}{8}, \tfrac{1}{8}, x)$$
$$\hat{f_2}(x) = g(-\tfrac{1}{8}, \tfrac{1}{8}, \tfrac{1}{8}, x) + g(-\tfrac{1}{4}, \tfrac{1}{4}, \tfrac{1}{4}, x)$$
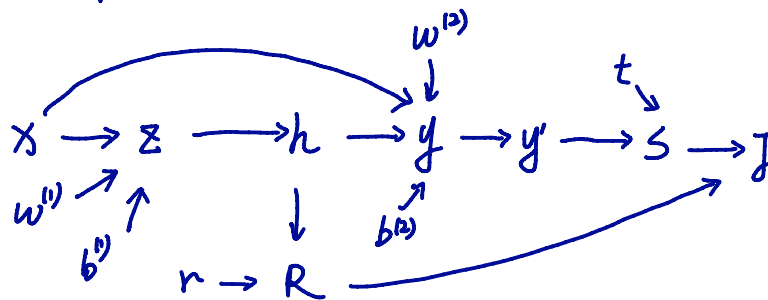$$\hat{f_3}(x) = g(-\tfrac{1}{8}, \tfrac{1}{8}, \tfrac{1}{8}, x) + g(-\tfrac{1}{4}, \tfrac{1}{4}, \tfrac{1}{4}, x) + g(-\tfrac{1}{2}, \tfrac{1}{2}, \tfrac{1}{2}, x)$$



# 2. Backprop.

## 2.1 Computational Graph

### 2.1.1.



### 2.1.2.

$$\bar{J} = \frac{dJ}{dJ} = 1$$

$$\bar{S} = \frac{dJ}{dS} = 1$$

$$\bar{y'} = \frac{dJ}{dy'} = \frac{dJ}{dS} \cdot \frac{dS}{dy'} = \bar{S} \cdot \mathbb{I}(t = k)$$

$$\bar{y} = \frac{dJ}{dy} = \frac{dJ}{dy'} \cdot \frac{dy'}{dy} = \bar{y'} \cdot \text{softmax}'(y)$$

$$\bar{R} = \frac{dJ}{dR} = 1$$

$$\bar{h} = \frac{dJ}{dh} = \frac{dJ}{dy} \cdot \frac{dy}{dh} + \frac{dJ}{dR} \cdot \frac{dR}{dh} = W^{(2)T} \bar{y} + \bar{R} r$$

$$\bar{z} = \frac{dJ}{dz} = \frac{dJ}{dh} \cdot \frac{dh}{dz} = \bar{z} = \begin{cases} \bar{h} & z > 0 \\ 0 & z \le 0 \end{cases}$$

$$\bar{x} = \frac{dJ}{dx} = \frac{dJ}{dz} \cdot \frac{dz}{dx} + \frac{dJ}{dy} \cdot \frac{dy}{dx}$$
$$= W^{(1)T} \bar{z} + \bar{y}$$

## 2.2. VJPs.

**2.2.1.**

$$f(x) = vv^Tx = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}[1\ 2\ 3]x = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 6 \\ 3 & 6 & 9 \end{bmatrix}x$$

$$\nabla f(x) = J = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 6 \\ 3 & 6 & 9 \end{bmatrix}$$

**2.2.2.**

$$J = v \cdot v^T \implies time: n^2 \text{ multiplication}$$
$$space: n^2$$

**2.2.3.**

$$J^Ty = [vv^T]^Ty = v^Tvy = v^Tyv \implies a = v^Ty, \quad J^Ty = a\cdot v \quad a \text{ is a scalar.}$$

$$\underbrace{}_{\substack{linear. \\ time + space}} + \underbrace{}_{\substack{linear \\ time + space}} \implies Linear\ time + space.$$

$$z = J^Ty \quad a = v^Ty = [1, 1, 1]\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} = 6$$

$$z = a\cdot v = 6v = \begin{bmatrix} 6 \\ 12 \\ 18 \end{bmatrix} \quad z^T = [6, 12, 18].$$

# 3. Linear Regression

**3.1.**

$$L = \frac{1}{n}\|x\hat{w} - t\|_2^2 = \frac{1}{n}(x\hat{w}-t)^T(x\hat{w}-t) = \frac{1}{n}(\hat{w}^Tx^T - t^T)(x\hat{w}-t) = \frac{1}{n}(\hat{w}^Tx^Tx\hat{w} - t^Tx\hat{w} - \hat{w}^Tx^Tt + t^Tt)$$

$$\frac{dL}{d\hat{w}} = \frac{1}{n}(2x^Tx\hat{w} - 2x^Tt) = \frac{2}{n}x^T(x\hat{w}-t)$$

**3.2.**

**3.2.1.**

$$\frac{dL}{d\hat{w}} = 0 \implies x^T(x\hat{w}-t) = 0 \quad x^Tx\hat{w} = x^Tt$$
$$\hat{w} = (x^Tx)^{-1}x^Tt$$

**3.2.2.** $t_i = w^{*T}x_i. \implies t = xw^*$

$$\hat{w} = (x^Tx)^{-1}x^Txw^* \implies \hat{w} = w^*$$

**3.3.**

**3.3.1.** $w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \quad w^Tx_1 = t. \quad 2w_1 + w_2 = 2 \implies$ infinity many $w_1, w_2$ satisfies.
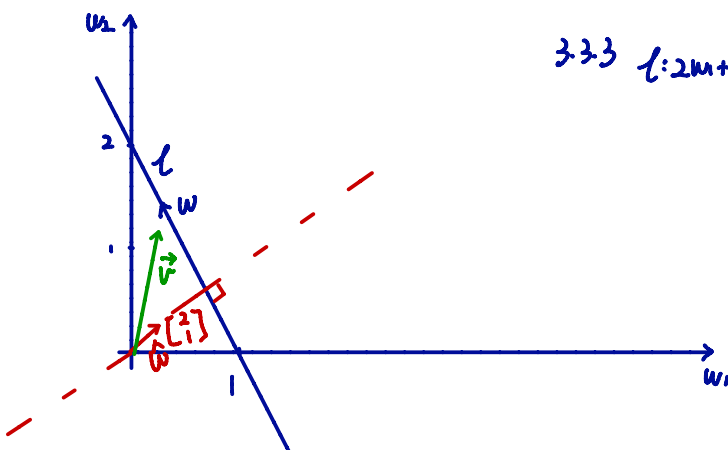
**3.3.2.**

$$\frac{dL}{d\hat{w}} = \frac{2}{n}x^T(x\hat{w}-t) = 2\cdot\begin{bmatrix} 2 \\ 1 \end{bmatrix}([2,1]\begin{bmatrix} 0 \\ 0 \end{bmatrix} - 2) = -4\begin{bmatrix} 2 \\ 1 \end{bmatrix}.$$

$$w_1 \leftarrow w_0 - d\frac{dL}{d\hat{w}} \quad w_1 = c\begin{bmatrix} 2 \\ 1 \end{bmatrix}.$$

$$\frac{dL}{d\hat{w}} = 2\cdot\begin{bmatrix} 2 \\ 1 \end{bmatrix}\cdot([2,1]\cdot c\begin{bmatrix} 2 \\ 1 \end{bmatrix} - 2) = k\begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

$$\implies w_2 = d\begin{bmatrix} 2 \\ 1 \end{bmatrix} \cdots \implies w_i. \text{ for } i \in \mathbb{Z}^+. \ w_i = a\begin{bmatrix} 2 \\ 1 \end{bmatrix}. \ a \in \mathbb{R}$$

$$\therefore \text{ all weight in direction of } \hat{w} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$



**3.3.3** $\ell : 2w_1 + w_2 = 2$ has direction: $w = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$

$$[1, -2]\cdot\begin{bmatrix} 2 \\ 1 \end{bmatrix} = 0$$

for any $\vec{v} = w - \hat{w}$

when $\vec{v} \perp w, \vec{v} = \hat{w}$, and has smallest

length / Euclidean norm by $\triangle$ enclosed by $\vec{v}, \hat{w}, w$

## 3.4.

### 3.4.1

$$\frac{dL}{dw} = \frac{2}{n} \sum_{1}^{n} x_i^T (w^T x_i - t_i)$$

$$\hat{w}_0 = 0$$

$$\frac{dL}{dw} = \frac{2}{n} \sum_{1}^{n} - t_i x_i^T$$

$$\hat{w}_1 = \hat{w}_0 - d\frac{dL}{dw} = c_1 x_1 + c_2 x_2 + \cdots + c_n x_n = X^T \cdot c$$

$$\frac{dL}{dw} = \frac{2}{n} \sum_{1}^{n} x_i^T (c^T X x_i - t_i) = X^T \cdot D \cdots$$

for each $\hat{w}_i$, $i \in \mathbb{Z}^+$, $\hat{w}_i = X^T \cdot C$, $C \in \mathbb{R}^{n \times 1}$

$$\hat{w} = X^T C \qquad X\hat{w} = t$$

$$\Rightarrow X X^T C = t \qquad c = (XX^T)^{-1} t$$

$$\hat{w} = X^T (XX^T)^{-1} t$$

### 3.4.2.

$$X \hat{w}_1 = t \qquad \hat{w}_1 = X^{-1} t \qquad \hat{w}_1^T = t^T (X^{-1})^T$$

$$\hat{w}^T = t^T [X^T (XX^T)^{-1}] = t^T [(XX^T)^{-1}]^T X$$

$$\hat{w}^T \hat{w} = t^T [(XX^T)^{-1}]^T X X^T (XX^T)^{-1} t = t^T (XX^T)^{-1} (XX^T)(XX^T)^{-1} t = t^T (XX^T)^{-1} t$$

$$\hat{w}_1^T \hat{w} = t^T (X^{-1})^T X^T (XX^T)^{-1} t = t^T (XX^T)^{-1} t$$

$$\therefore (\hat{w} - \hat{w}_1)^T \hat{w} = \hat{w}^T \hat{w} - \hat{w}_1^T \hat{w} = 0$$

$\therefore \hat{w}$ has smallest distance / euclidean norm among all possible $w$ $\Rightarrow \hat{w}_1 \perp (\hat{w} - \hat{w}_1)$

```
# to be implemented; fill in the derived solution for the underparameterized (d<n) and overparameterized (d>n) problem

def fit_poly(X, d, t):
    X_expand = poly_expand(X, d=d, poly_type = poly_type)
    if d > n:
        W = np.dot(np.dot(np.transpose(X_expand), np.linalg.inv(np.dot(X_expand, np.transpose(X_expand)))), t)
    else:
        W = np.dot(np.dot(np.linalg.inv(np.dot(np.transpose(X_expand), X_expand)), np.transpose(X_expand)), t)
    return W
```

Higher degree polynomial / Overparameterization does not always leads to overfitting.

degree 70, 100, 150, 200 showes better generization and smaller error than degree 10-50.