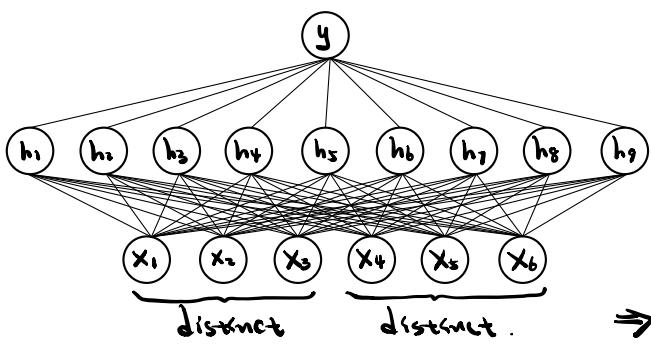


1. Hard-Coding Networks.

1.1 Verify Element in List

$$W^1 = \begin{bmatrix} -2 & 0 & 0 & 2 \\ 0 & -2 & 0 & 2 \\ 0 & 0 & -2 & 2 \end{bmatrix} \quad b^1 = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \quad W^2 = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \quad b^2 = -2$$

1.2 Verify Permutation



⇒ Just like the neural network in Section 1.1, it has 4 inputs which compares the last number to the other three numbers and output a 1 to hidden layer when there is a match.

⇒ Same idea applies here, this neural network can compare x_4, x_5 and x_6 with other 3 numbers x_1, x_2 and x_3 , respectively. By applying proper weights and bias terms, and same activation function, the outputs of hidden layer could be: (output 0 if no match).

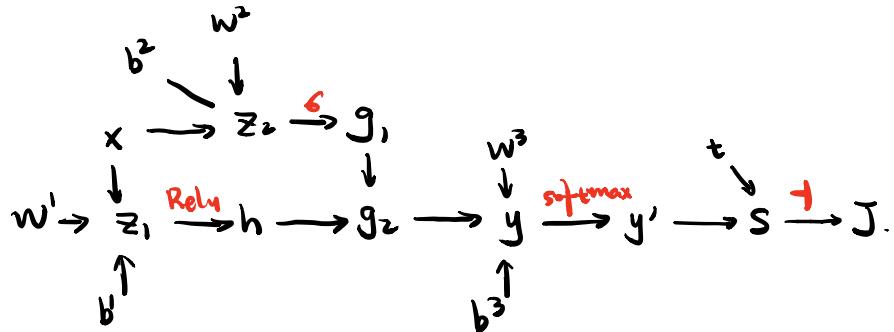
- { x_4 vs. x_{1-3} : output to h_{1-3} ; has a 1 if x_4 matches either one from x_{1-3}
- x_5 vs. x_{1-3} : output to h_{4-6} ; has a 1 if x_5 matches either one from x_{1-3}
- x_6 vs. x_{1-3} : output to h_{7-9} ; has a 1 if x_6 matches either one from x_{1-3}

If x_{4-6} are a permutation of x_{1-3} , the sum of all hidden units will be 3 and output a 1 to the output layer by applying proper parameters.

2. Backpropagation.

2.1

{2.1.1 computational graph.}



{2.1.2 Backward Pass}

$$\mathbf{z}_1 = \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}$$

$$\mathbf{h} = \text{ReLU}(\mathbf{z}_1)$$

$$\mathbf{z}_2 = \mathbf{W}^{(2)}\mathbf{h} + \mathbf{b}^{(2)}$$

$$g_1 = \sigma(\mathbf{z}_2)$$

$$g_2 = \mathbf{h} \odot g_1$$

$$\mathbf{y} = \mathbf{W}^{(3)}\mathbf{g}_2 + \mathbf{b}^{(3)}$$

$$\mathbf{y}' = \text{softmax}(\mathbf{y})$$

$$\mathcal{S} = \sum_{k=1}^N \mathbb{I}(t=k) \log(y'_k)$$

$$\mathcal{J} = -\mathcal{S}$$

$$\bar{\mathcal{J}} = \frac{\partial \mathcal{J}}{\partial \mathcal{J}} = 1$$

$$\bar{s} = \bar{\mathcal{J}} \frac{\partial \mathcal{J}}{\partial s} = -\bar{\mathcal{J}}$$

$$\bar{y}' = \bar{s} \frac{\partial s}{\partial y'} = (-1) \sum_{k=1}^N \mathbb{I}(t=k) \frac{1}{y'_k} = \sum_{k=1}^N \mathbb{I}(t=k) \frac{-1}{y'_k}$$

$$\bar{y} = \bar{y}' \frac{\partial y'}{\partial y} = \bar{y}' \odot \text{softmax}'(y)$$

$$= \begin{cases} \bar{y}' [\text{softmax}(y_k) \times (1 - \text{softmax}(y_j))] & k=j \\ \bar{y}' [-\text{softmax}(y_k) \times \text{softmax}(y_j)] & k \neq j \end{cases}$$

$$\bar{g}_2 = \bar{y} \frac{\partial y}{\partial g_2} = W^3 \bar{y}$$

$$\bar{g}_1 = \bar{g}_2 \frac{\partial g_2}{\partial g_1} = \bar{g}_2 \odot h$$

$$\bar{h} = \bar{g}_2 \frac{\partial g_2}{\partial h} = \bar{g}_2 \odot g_1$$

$$\bar{z}_2 = \bar{g}_1 \frac{\partial g_1}{\partial z_2} = \bar{g}_1 \odot \sigma'(z_2)$$

$$= \bar{g}_1 \odot [\sigma(z_2) \cdot (1 - \sigma(z_2))]$$

$$\bar{z}_1 = \bar{h} \frac{\partial h}{\partial z_1} = \bar{h} \odot \text{ReLU}(z_1)$$

$$= \bar{h} \odot \begin{cases} 1 & \text{when } z_1 > 0 \\ 0 & \text{when } z_1 \leq 0 \end{cases} \quad \text{for } i \in \{1, \dots, M\}$$

$$\bar{x} = \bar{z}_1 \frac{\partial \bar{z}_1}{\partial x} + \bar{z}_2 \frac{\partial \bar{z}_2}{\partial x} = w^T \bar{z}_1 + w^T \bar{z}_2$$

2.2 Automatic Differentiation.

2.2.1 Compute Hessian. $L(x) = x^T V V^T x$

$$\Rightarrow \nabla L(x) = 2V V^T x \Rightarrow H(x) = 2V V^T$$

$$= 2 \times \begin{bmatrix} 4 \\ 2 \\ 3 \end{bmatrix} \begin{bmatrix} 4 & 2 & 3 \end{bmatrix} = \begin{bmatrix} 32 & 16 & 24 \\ 16 & 8 & 12 \\ 24 & 12 & 18 \end{bmatrix}$$

2.2.2 Computation cost

\Rightarrow number of scalar multiplication : $O(n^2)$

\Rightarrow memory cost : $O(n^2)$

2.3 Vector-Hessian Products.

$$\Rightarrow \text{Reverse Mode: } M = V^T y = [1 \ 2 \ 3] \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = 1+2+3 = 6$$

$$z = V M = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \times 6 = \begin{bmatrix} 6 \\ 12 \\ 18 \end{bmatrix}$$

\Rightarrow time and memory cost = $O(n)$

$$\Rightarrow \text{Forward Mode: } H = V V^T = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} [1 \ 2 \ 3] = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 6 \\ 3 & 6 & 9 \end{bmatrix}$$

$$z = Hy = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 6 \\ 3 & 6 & 9 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 6 \\ 12 \\ 18 \end{bmatrix}$$

\Rightarrow time and memory cost = $O(n^2)$

2.4 Trade-off of Reverse- and Forward-mode Autodiff

$$\Rightarrow \text{Reverse mode: } Z = V(V^T(y_1 y_2^T)) \\ = [n \times 1] ([1 \times n] ([n \times 1] [1 \times m]))$$

\Rightarrow computation cost = $3nm$.

$$\Rightarrow \text{Forward mode: } Z = VV^T y_1 y_2^T \\ = [n \times 1] [1 \times n] [n \times 1] [1 \times m]$$

\Rightarrow computation cost = $2n^2 + nm$.

$\left\{ \begin{array}{l} \text{wide } (m \gg n): m \text{ dominates } n \Rightarrow 2n^2 + nm < 3nm \Rightarrow \text{forward mode is better.} \\ \text{tall } (n \gg m): n \text{ dominate } m \Rightarrow 2n^2 + m > 3nm \Rightarrow \text{Backward mode is better} \end{array} \right.$

3. Linear Regression

$$\min_{\hat{w}} \frac{1}{n} \sum_{i=1}^n (\hat{w}^T x_i - t_i)^2 = \min_{\hat{w}} \frac{1}{n} \|X\hat{w} - t\|_2^2.$$

3.1 Deriving the Gradient.

$$\begin{aligned} L &= \frac{1}{n} \|X\hat{w} - t\|_2^2 = \frac{1}{n} (X\hat{w} - t)^T (X\hat{w} - t) \\ &= \frac{1}{n} (\hat{w}^T X^T - t^T) (X\hat{w} - t) \\ &= \frac{1}{n} (\hat{w}^T X^T X \hat{w} - \hat{w}^T X^T t - t^T X \hat{w} + t^T t) \end{aligned}$$

$$\Rightarrow \frac{\partial L}{\partial \hat{w}} = \frac{1}{n} (2X^T X \hat{w} - X^T t - X^T t + 0) \\ = \underbrace{\frac{2}{n} (X^T X \hat{w} - X^T t)}_{= 0}.$$

3.2 Underparameterized model. ($d < n$)

$$\text{set } \frac{\partial L}{\partial \hat{w}} = 0 \Rightarrow \frac{2}{n} (X^T X \hat{w} - X^T t) = 0.$$

$$X^T X \hat{w} - X^T t = 0.$$

$$X^T X \hat{w} = X^T t.$$

$$\hat{w} = \underbrace{(X^T X)^{-1} X^T t}_{= 0}.$$

$X^T X$ is invertible when $n > d \Rightarrow$ solution is unique.

$$0 \quad d \quad d \quad d$$

3.3 reparameterize mo ($L > n$)

3.3.1 Let $\hat{w} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$.

$$\Rightarrow \hat{w}^T x_i = t_i$$

Substitute vectors: $[w_1 \ w_2] \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 3 \Rightarrow \underline{w_2 = 3 - w_1}$

\Rightarrow therefore, there are infinite number of weights on this line.

3.3.2 from 3.1 $\Rightarrow \frac{\partial L}{\partial \hat{w}} = \frac{2}{n} (X^T X \hat{w} - X^T t)$.
 $= \frac{2}{n} X^T (X \hat{w} - t)$.

initialization $\hat{w} = 0 \Rightarrow \frac{\partial L}{\partial \hat{w}_0} = \frac{-2}{n} X^T t$

first update: $\hat{w}_1 \leftarrow \hat{w}_0 - \alpha \frac{\partial L}{\partial \hat{w}_0} = 0 + \frac{2\alpha}{n} X^T t = \frac{2\alpha}{n} X^T t$.

$\Rightarrow \frac{2\alpha}{n} X^T t$ can be written as $\underline{X^T C_1}$ for a vector of constants
with dimension $C_1 \in \mathbb{R}^{n \times 1}$.

second update: $\frac{\partial L}{\partial \hat{w}_1} = \frac{2}{n} X^T (X \hat{w}_1 - t)$
 $= \frac{2}{n} X^T (X(X^T C_1) - t)$.

$$\hat{w}_2 \leftarrow \hat{w}_1 - \alpha \frac{\partial L}{\partial \hat{w}_1} = X^T C_1 - \frac{2\alpha}{n} X^T (X(X^T C_1) - t)
= X^T \underbrace{(C_1 - \frac{2\alpha}{n} (X(X^T C_1) - t))}_{\text{still a vector of constants in } \mathbb{R}^{n \times 1}}$$

therefore, \hat{w}_2 can be written as $X^T C_2$ for $C_2 \in \mathbb{R}^{n \times 1}$.

Conclusion: as we can see from the two weight updates

$$\Rightarrow \text{for each } \hat{w}_i, \hat{w}_i = X^T C \text{ for } X \in \mathbb{R}^{n \times d}, C \in \mathbb{R}^{d \times 1}$$

therefore, from $\hat{w} = X^T c$ and $X \hat{w} = t$

①

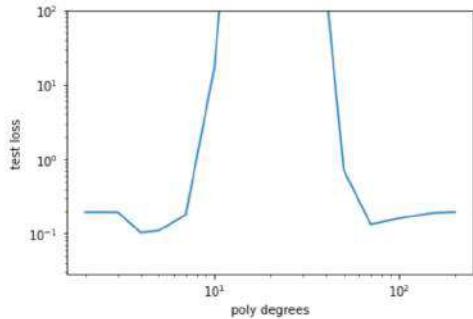
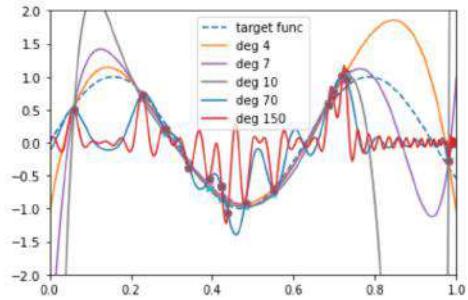
②

Substitute ① to ② : $X X^T c = t \Rightarrow c = (X X^T)^{-1} t$
 $\Rightarrow \hat{w} = X^T c = \underbrace{X^T (X X^T)^{-1} t}_{\text{Part 3.3.2}}$.

3.3.3

► # 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)
    n = X.shape[0]
    if d > n:
        ## W = ... (Your solution for Part 3.3.2)
        W = np.dot(np.dot(np.transpose(X_expand),np.linalg.inv(np.dot(X_expand,np.transpose(X_expand)))),t)
    else:
        ## W = ... (Your solution for Part 3.2)
        W = np.dot(np.dot(np.linalg.inv(np.dot(np.transpose(X_expand),X_expand)),np.transpose(X_expand)),t)
    return W
```



Overparameterization (high degree polynomial) does not always lead to overfitting. By observing the second graph, polynomial degree around 70 or bigger shows better generalization in test data than degree 10 - 40 (underparameterization)