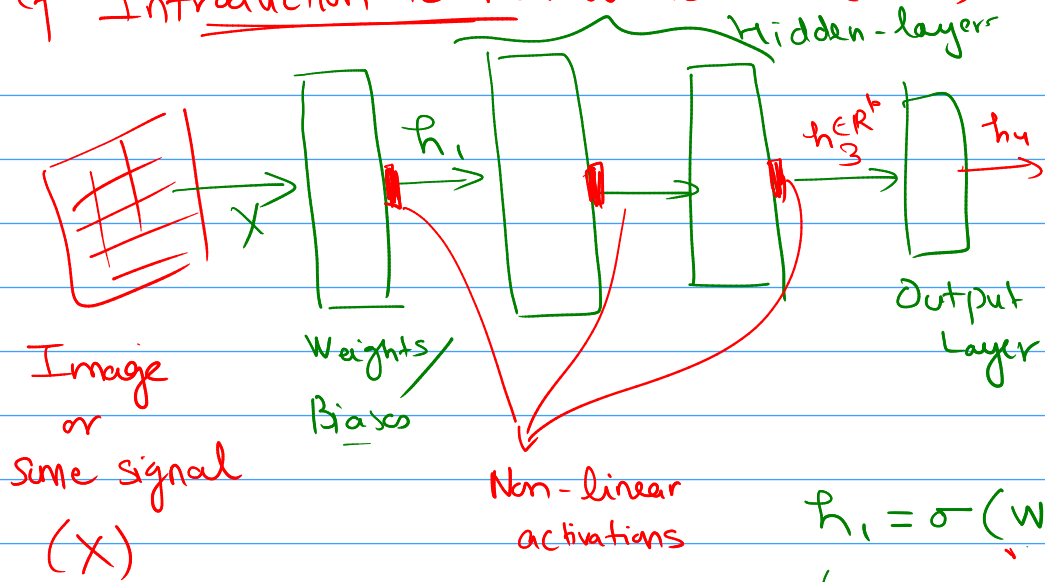


* Brief Introduction to Neural Networks (NNs)



$$\tilde{X} = \begin{bmatrix} X \\ 1 \end{bmatrix}$$

$W \in \mathbb{R}^{m \times (n+1)}$

$$W \in \mathbb{R}^{m \times n}$$

$$b \in \mathbb{R}^m$$

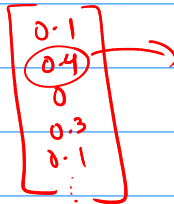
$$h_1 = \sigma(WX + b)$$

\mathbb{R}^m \mathbb{R}^n

Output layer → usually of size $1 \times p$ (Regression)

→ K-fold classification problem ($K \times p$)

K classes



$$\sum p_i = 1$$

$$u_h = W h_3$$

logits

Soft-max activation

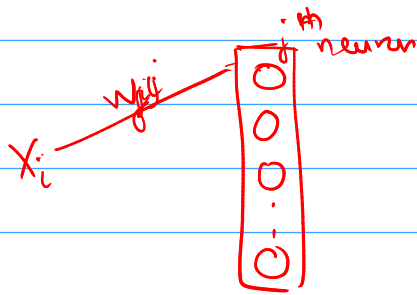
$$y_i = \frac{\exp(u_h(i))}{\sum_j \exp(u_h(j))}$$

$(h_i = \max_{i \in [p]} \{u_h(i)\})$ → then take softmax

$$Y = f_{\theta}(X)$$

→ θ : weights of the NN

Theorem: Any piecewise continuous function can be approximated arbitrarily close using a four-layer NN with ReLU activations



$$\{x_i, y_i\}_{i=1}^N$$

$W \rightarrow$ weights of a NN

* How do we train a NN?

\hookrightarrow Loss functions
 $\mathcal{L}(W; \xi)$

Regression $\rightarrow \frac{(\hat{y}_i - y_i)^2}{2}$

Cross-entropy loss.

Classification

$$-\sum_{i=1}^k t_i \log p_i$$

\hookrightarrow minimize this loss function

Gradient-descent algorithm:

$$\min_{x \in \mathbb{R}^n} f(x)$$

$$\boxed{x^{(k+1)} = x^{(k)} - \eta_k \nabla_x f(x^{(k)})} \rightarrow \text{Gradient descent}$$

Stochastic Gradient descent (SGD):

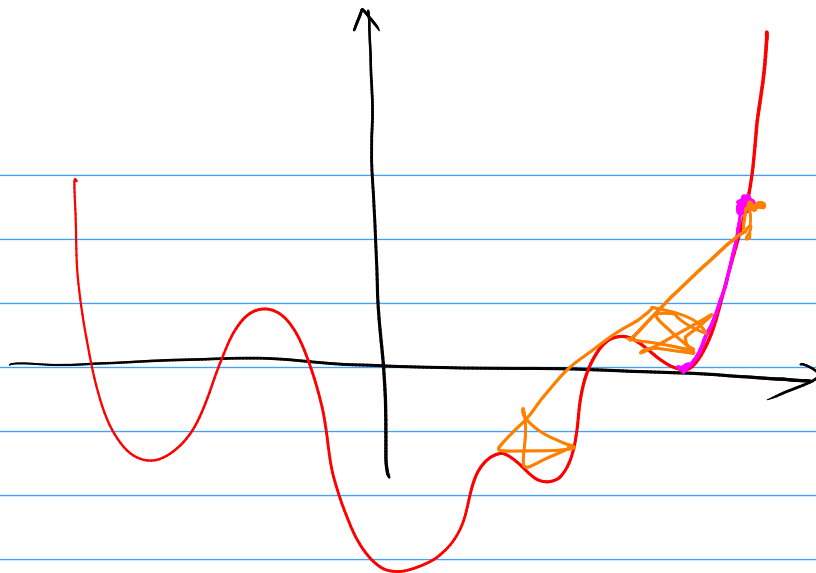
$$W^{(k+1)} = W^{(k)} - \eta_k \nabla_{W_k} L(W^{(k)}, \xi_k)$$

\hookrightarrow Data points sampled in k^{th} iteration

Full-batch Gradient descent (GD):

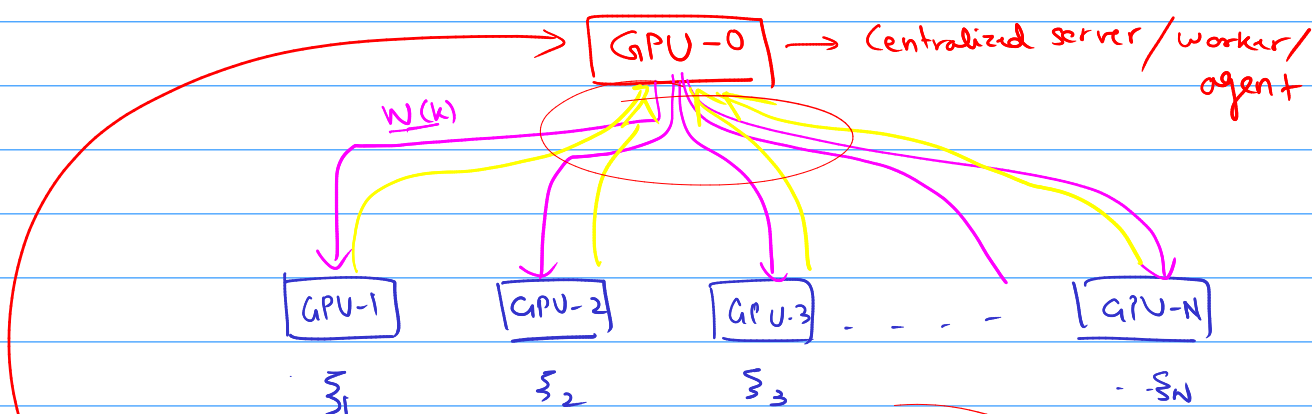
$$W^{(k+1)} = W^{(k)} - \eta_k \frac{1}{N} \sum_{i=1}^N \nabla_{W_k} l_{w_i}(W^{(k)}, x_i)$$

\hookrightarrow Too-slow (because of gradient computation on all data points)



* How do we work with large scale data?

Parameter - server approach



- Privacy concerns
- Limited memory

$$w(k+1) = w(k) - \eta \sum_{i=1}^N \frac{1}{N_i} \nabla_w L(w(k), \xi_i)$$

- Single point of failure.
- Communication bandwidth requirement scales with the number of servers

* (Baidu) Ring All Reduce algorithm

↳ communication bandwidth requirement is constant in the number of agents.

↳ There is no centralized server.

Scatter Reduce

All Gather

Example: 3 workers/agents

$(W^{(k)})$

1st iteration

$a_0 \mid b_0 \mid c_0$

1st agent

$a_1 \mid b_1 \mid c_1$

2nd agent

$a_2 \mid b_2 \mid c_2$

3rd agent

2nd iteration

$a_0 \mid b_0 \mid c_0 + c_2$

$a_1 + a_0 \mid b_1 \mid c_1$

$a_2 \mid b_1 + b_2 \mid c_2$

$$2(N-1) \times \frac{K}{N} = \Omega(1)$$

1st

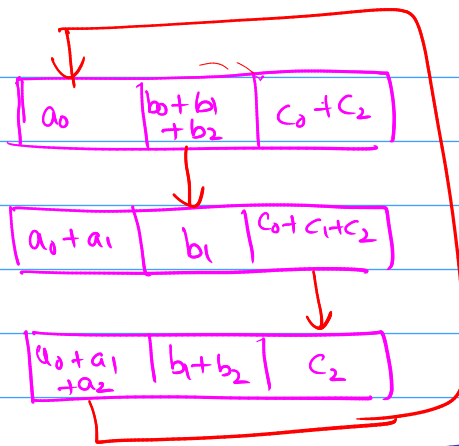
$a_0 \mid b_0 + b_1 + b_2 \mid c_0 + c_2$

$a_0 + a_1 \mid b_1 \mid c_0 + c_1 + c_2$

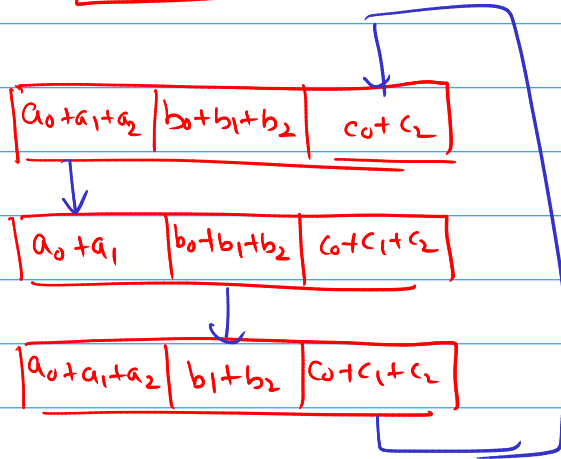
$a_0 + a_1 + a_2 \mid b_1 + b_2 \mid c_2$

All Gather

1st iteration



2nd iteration



Result on SGD:

* Under the assumptions:

(i) The loss function $F(w; \xi)$ is L -smooth in terms of w .

(ii) Stochastic Gradient is unbiased and has bounded variance σ^2 .

$$\mathbb{E}_{\xi \sim D} [\nabla F(w^{(k)}; \xi)] = \nabla \mathbb{E}_{\xi \sim D} [F(w^{(k)}; \xi)]$$

Then, $\mu = \mathcal{O}\left(\frac{1}{\sqrt{T}}\right) = \nabla f(w^{(k)})$

$$\frac{1}{T} \sum_{k=0}^{T-1} \mathbb{E} \|\nabla f(w^{(k)})\|^2 = \mathcal{O}\left(\frac{\sigma}{\sqrt{T}}\right)$$

ϵ