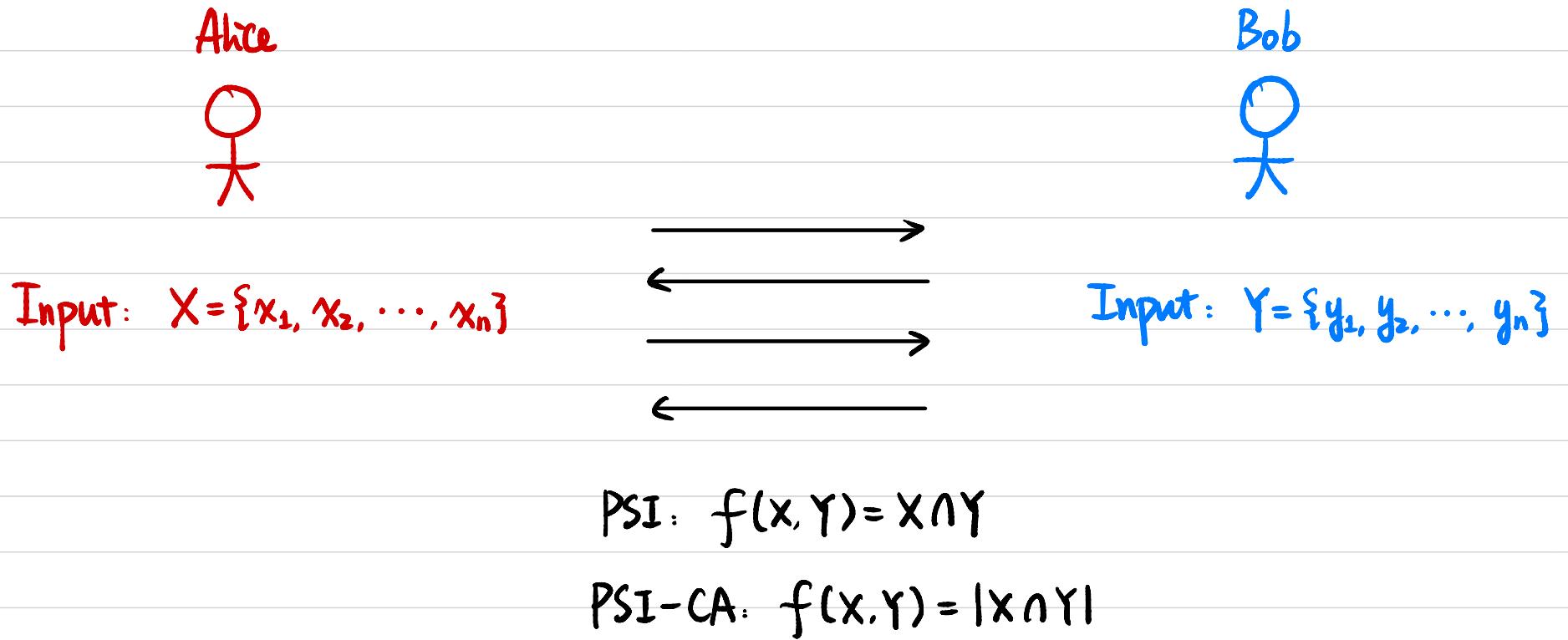


CSCI 1515 Applied Cryptography

This Lecture:

- Private Set Intersection
- Privacy-Preserving Machine Learning

Private Set Intersection (PSI)



Applications:

- Password Breach Alert (Chrome, Edge, Firefox, iOS Keychain, ...)
- Ads Conversion Measurement (Google)
- Privacy-Preserving Inventory Matching (J.P. Morgan)
- Private Database Joins

Private Set Intersection (PSI)

Alice



Input: $X = \{x_1, x_2, \dots, x_n\}$

$H(x_1), \dots, H(x_n)$

Bob



Input: $Y = \{y_1, y_2, \dots, y_n\}$

$H(y_1), \dots, H(y_n)$



$X \cap Y$

Is it (semi-honest) secure?

Is it possible to achieve 2PC / MPC with 1 round of communication?

DDH-based PSI Cyclic group G of order q with generator g , where DDH holds.
 $H: \{0,1\}^* \rightarrow G$ (modeled as Random Oracle)

Alice



Bob



Input: $X = \{x_1, x_2, \dots, x_n\}$

$$a \leftarrow \mathbb{Z}_q$$

$$\leftarrow H(Y)^b := \{H(y_1)^b, \dots, H(y_n)^b\}$$

$$\overbrace{H(X)^a, H(Y)^{a \cdot b}}$$

$$b \leftarrow \mathbb{Z}_q$$

Input: $Y = \{y_1, y_2, \dots, y_n\}$

$$H(X)^{a \cdot b} \cap H(Y)^{a \cdot b}$$

Is it (semi-honest) secure?

$$\downarrow \\ X \cap Y$$

PSI-CA ?

PSI-CA: $f(x, y) = |x \cap y|$

Alice



Bob



Input: $X = \{x_1, x_2, \dots, x_n\}$

$$a \leftarrow \mathbb{Z}_q$$

$$\underbrace{H(Y)^b := \{H(y_1)^b, \dots, H(y_n)^b\}}$$

Input: $Y = \{y_1, y_2, \dots, y_n\}$

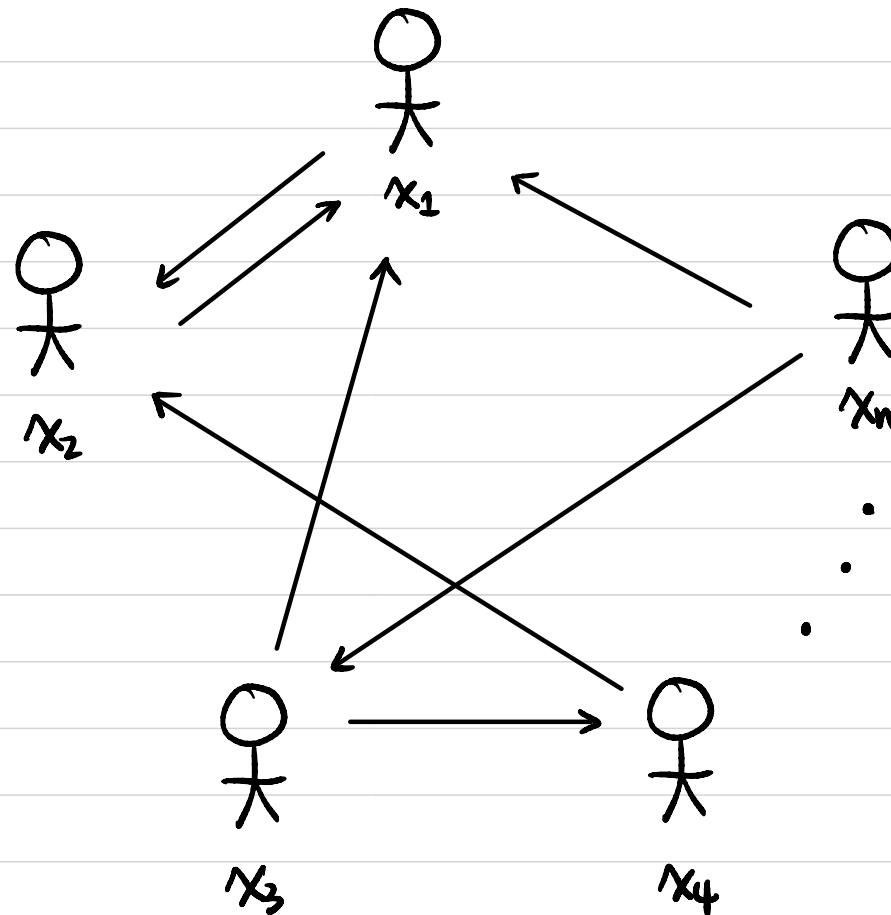
$$b \leftarrow \mathbb{Z}_q$$

$$\overbrace{H(X)^a, H(Y)^{a \cdot b}}$$

$$H(X)^{a \cdot b} \cap H(Y)^{a \cdot b}$$

$$\downarrow \\ X \cap Y$$

Privacy-Preserving Machine Learning (PPML)

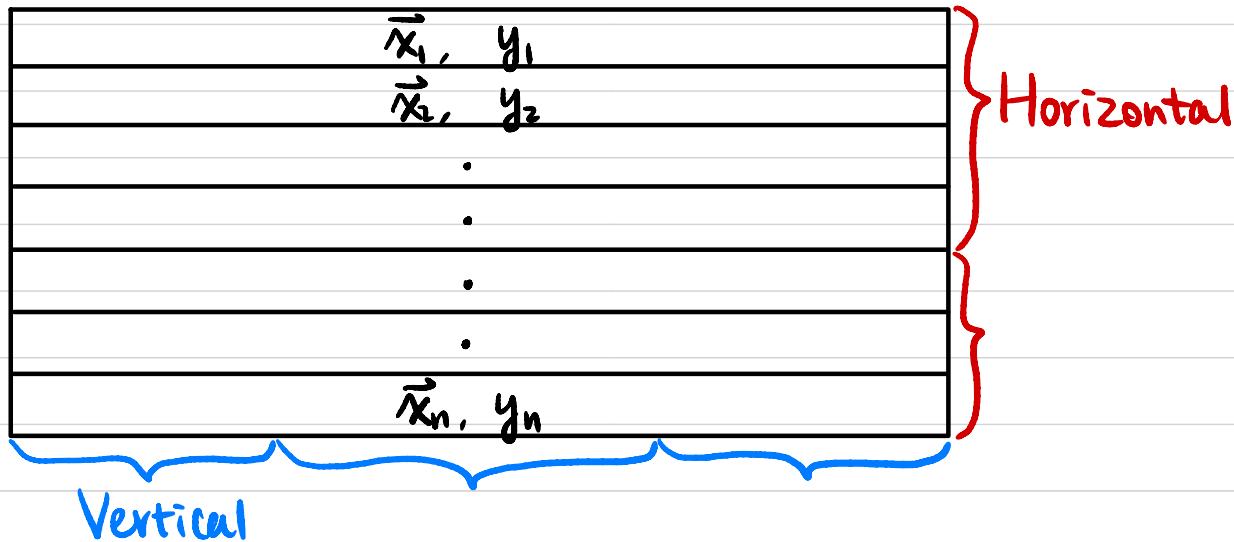


↓
ML model

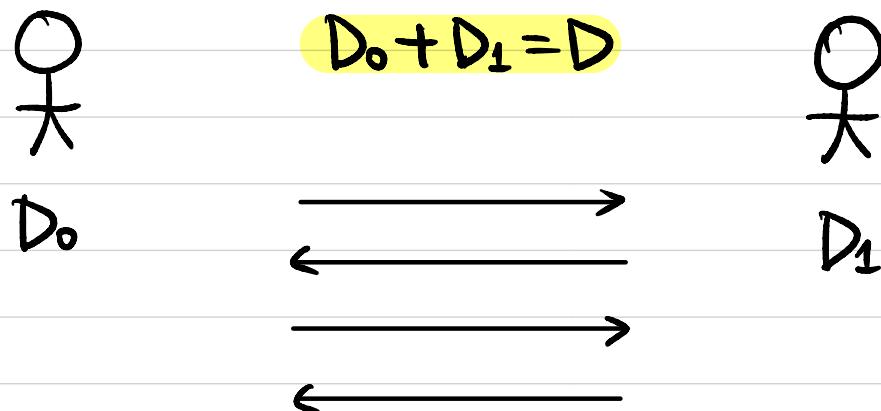
↓
ML inference

Privacy-Preserving Machine Learning (PPML)

Horizontal / Vertical Data Partitioning



PPML Framework



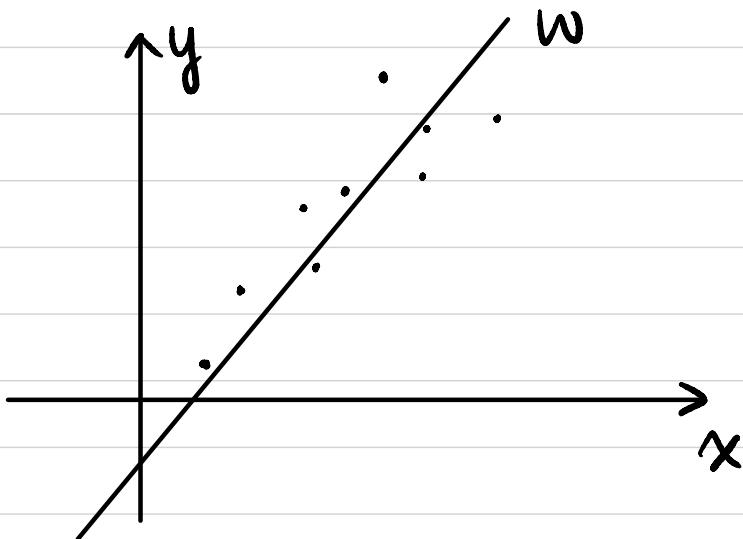
Machine Learning Background

Linear Regression

Data Points (\vec{x}, y)

ML Model: Coefficient vector \vec{w}

$$g(\vec{x}) = \langle \vec{x}, \vec{w} \rangle$$



For each data point (\vec{x}_i, y_i) ,
Define loss function

$$L_i(\vec{w}) := \frac{1}{2} (\langle \vec{x}_i, \vec{w} \rangle - y_i)^2$$

$$\text{Total loss: } L(\vec{w}) := \frac{1}{N} \sum_{i \in [N]} L_i(\vec{w})$$

Goal: Find \vec{w} that minimizes $L(\vec{w})$.

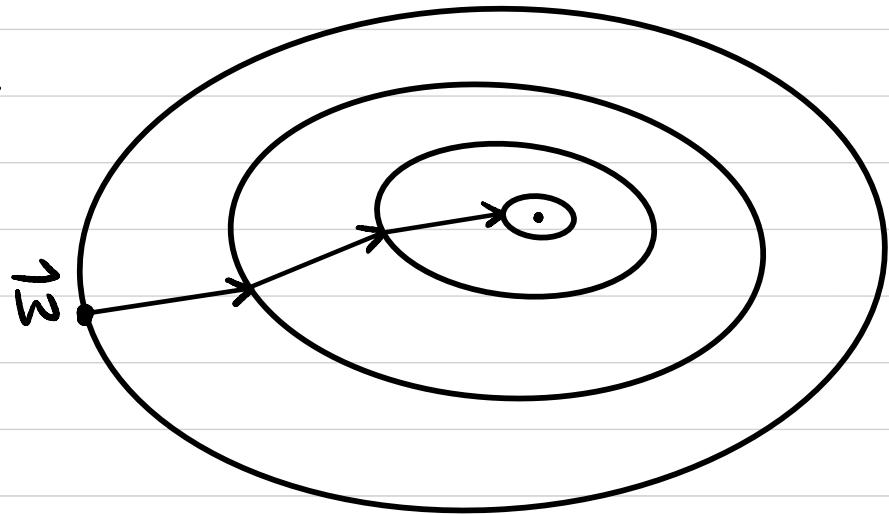
Machine Learning Background

Stochastic Gradient Descent (SGD)

- \vec{w} initialized with arbitrary value
- Given a data point (\vec{x}_i, y_i) :

$$\vec{w} \leftarrow \vec{w} - \eta \cdot \nabla L_i(\vec{w})$$

↑ ↑
learning rate gradient

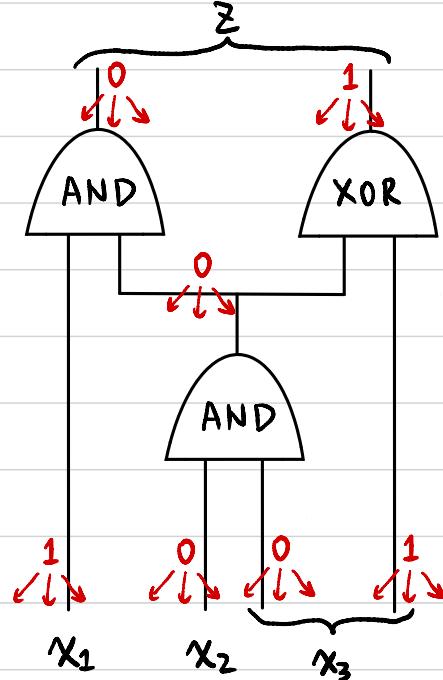


$$\vec{w} \leftarrow \vec{w} - \eta \cdot (\underbrace{\langle \vec{x}_i, \vec{w} \rangle - y_i}_{y_i^* = \langle \vec{x}_i, \vec{w} \rangle} \cdot \vec{x}_i)$$

$y_i^* = \langle \vec{x}_i, \vec{w} \rangle$ (forward propagation)

backward propagation

MPC for any function with $t \leq n-1$ (GMW)



Each party P_i holds a random share $V_i^w \in \{0, 1\}$ s.t. $\bigoplus_{i=1}^n V_i^w = v^w$

Inputs:

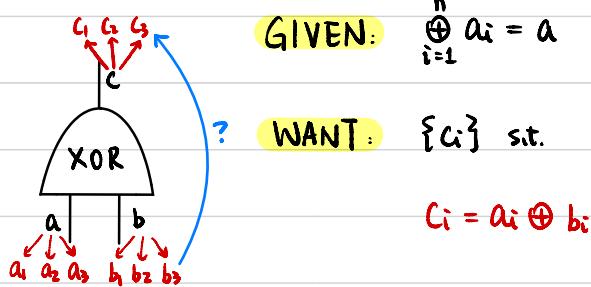
For each input wire w :

If it's from party P_k with input value $v^w \in \{0, 1\}$.

P_k randomly samples $V_i^w \xleftarrow{\$} \{0, 1\}$ s.t. $\bigoplus_{i=1}^n V_i^w = v^w$

Sends V_i^w to party P_i .

XOR gates:



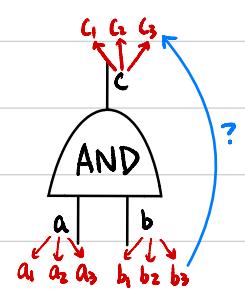
GIVEN: $\bigoplus_{i=1}^n a_i = a$

$\bigoplus_{i=1}^n b_i = b$

WANT: $\{c_i\}$ s.t. $\bigoplus_{i=1}^n c_i = c = a \oplus b$

$$c_i = a_i \oplus b_i$$

AND gates:



GIVEN: $\bigoplus_{i=1}^n a_i = a$

$\bigoplus_{i=1}^n b_i = b$

WANT: $\{c_i\}$ s.t. $\bigoplus_{i=1}^n c_i = c = a \cdot b$

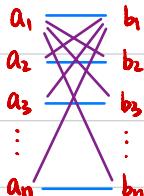
$$c_i = a_i \cdot b_i \oplus \sum_{j \neq i} r_{ij}^{(w)} \oplus \sum_{j \neq i} r_{ji}^{(w)}$$

$$a \cdot b = \left(\sum_{i=1}^n a_i \right) \cdot \left(\sum_{i=1}^n b_i \right) \pmod{2}$$

$$= \left(\sum_{i=1}^n a_i \cdot b_i \right) + \left(\sum_{i+j} a_i \cdot b_j \right) \pmod{2}$$

P_i locally

$r_{ij}^{(w)}$ Reshare $r_{ij}^{(w)}$



Outputs:

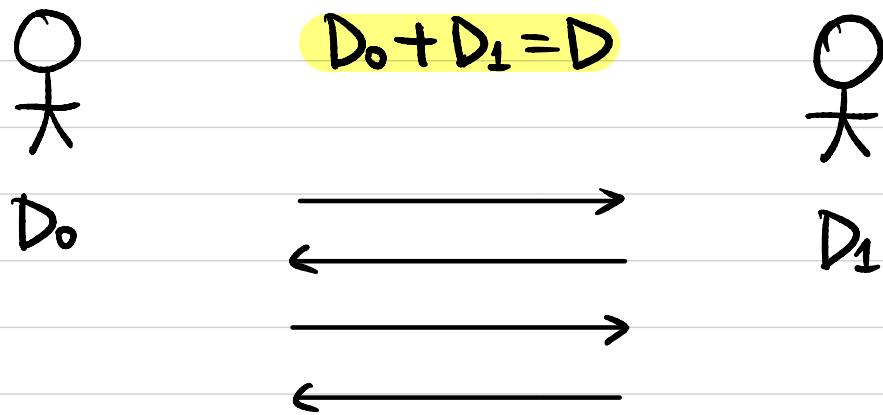
For each output wire w :

Each party P_i holds a random share $V_i^w \in \{0, 1\}$

Sends V_i^w to all parties

Each party computes the value $v^w = \bigoplus_{i=1}^n V_i^w$

PPML for Linear Regression



Invariant: Two parties hold secret shares of \vec{w} .

Initialization: Randomly sample \vec{w}_0, \vec{w}_1 respectively.

SGD: $\vec{w} \leftarrow \vec{w} - \eta \cdot (\langle \vec{x}_i, \vec{w} \rangle - y_i) \cdot \vec{x}_i$

How to get a secret share of the updated \vec{w} ?

Machine Learning Background

Logistic Regression

Data Points (\vec{x}, y)

ML Model: Coefficient vector \vec{w}

$$g(\vec{x}) = f(\langle \vec{x}, \vec{w} \rangle)$$

↑
activation function

Sigmoid $f(u) = \frac{1}{1+e^{-u}}$



For each data point (\vec{x}_i, y_i) ,

Define loss function $L_i(\vec{w}) := -y_i \cdot \log y_i^* - (1-y_i) \cdot \log (1-y_i^*)$

Total loss: $L(\vec{w}) := \frac{1}{N} \sum_{i \in [N]} L_i(\vec{w})$

$y_i^* := f(\langle \vec{x}_i, \vec{w} \rangle)$

Goal: Find \vec{w} that minimizes $L(\vec{w})$.

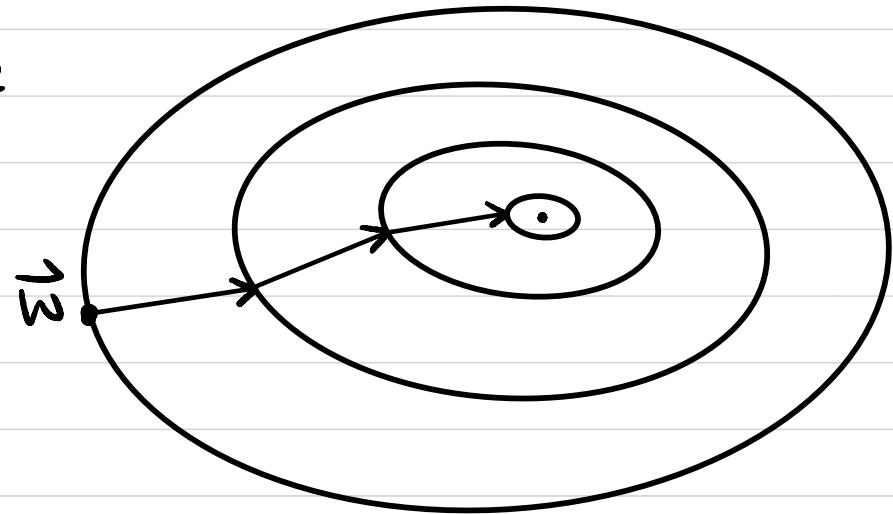
Machine Learning Background

Stochastic Gradient Descent (SGD)

- \vec{w} initialized with arbitrary value
- Given a data point (\vec{x}_i, y_i) :

$$\vec{w} \leftarrow \vec{w} - \eta \cdot \nabla L_i(\vec{w})$$

↑ ↑
learning rate gradient

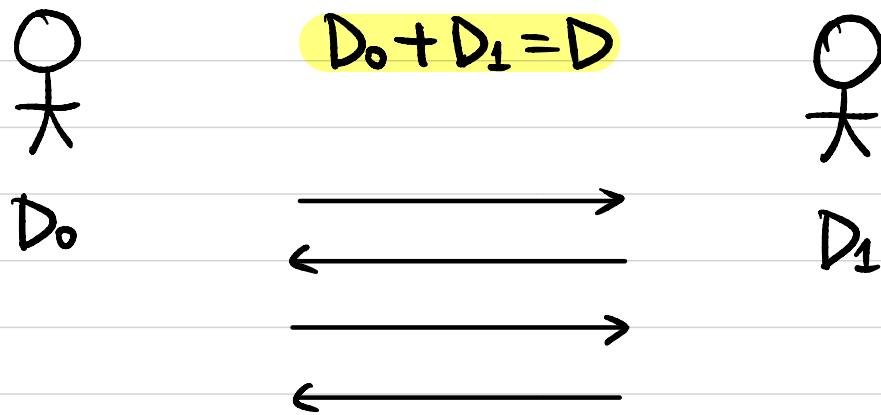


$$\vec{w} \leftarrow \vec{w} - \eta \cdot (f(\underbrace{\langle \vec{x}_i, \vec{w} \rangle}_{y_i^*}) - y_i) \cdot \vec{x}_i$$

$y_i^* = \langle \vec{x}_i, \vec{w} \rangle$ (forward propagation)

backward propagation

PPML for Logistic Regression



Invariant: Two parties hold secret shares of \vec{w} .

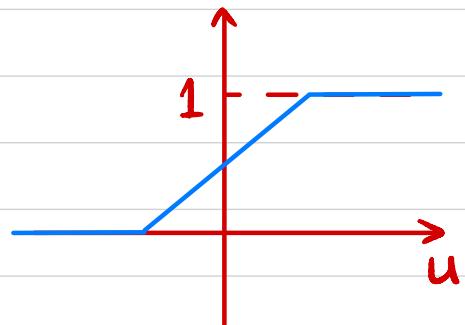
Initialization: Randomly sample \vec{w}_0, \vec{w}_1 respectively.

SGD: $\vec{w} \leftarrow \vec{w} - \eta \cdot (f(\langle \vec{x}_i, \vec{w} \rangle) - y_i) \cdot \vec{x}_i$

How to get a secret share of the updated \vec{w} ?

PPML for Logistic Regression

Approximating $f(u)$



Piecewise polynomial $f(u) = \begin{cases} 0 & \text{if } u < -\frac{1}{2} \\ u + \frac{1}{2} & \text{if } u \in [-\frac{1}{2}, \frac{1}{2}] \\ 1 & \text{if } u > \frac{1}{2} \end{cases}$

$$b_1 := \begin{cases} 1 & \text{if } u < -\frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$

$$b_2 := \begin{cases} 1 & \text{if } u > \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$$

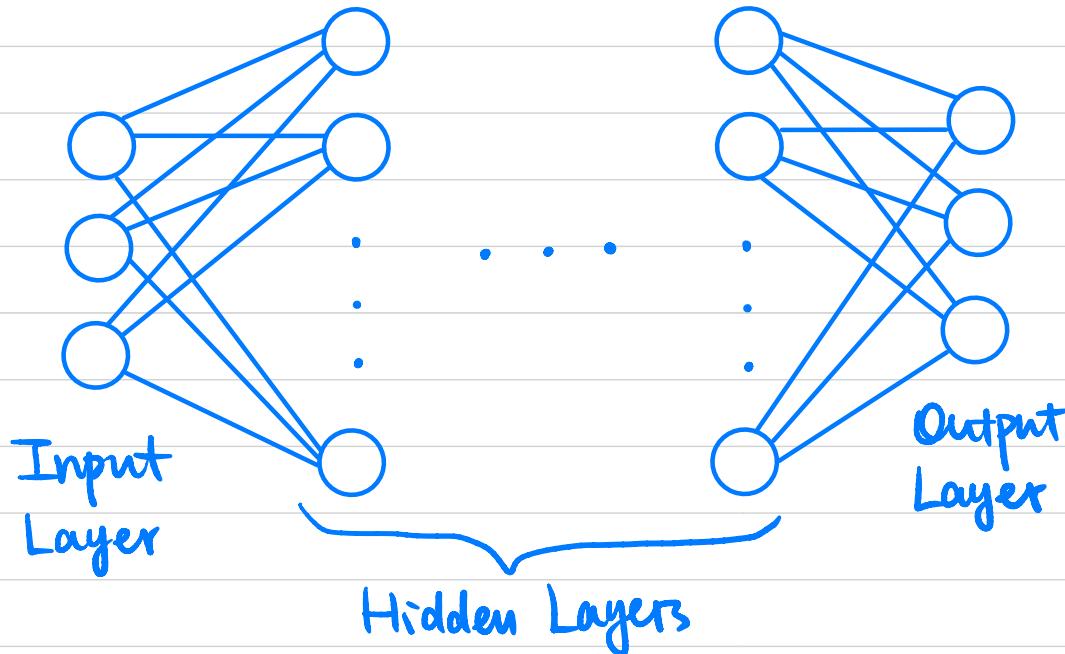
$$f(u) = ?$$

Machine Learning Background

Neural Network

Data Points (\vec{x}, y)

ML Model: coefficient vector \vec{w}



Each node in hidden layers: linear function + activation function

$$\text{Total loss: } L(\vec{w}) := \frac{1}{N} \sum_{i \in [N]} L_i(\vec{w})$$

Goal: Find \vec{w} that minimizes $L(\vec{w})$.