

The Sigmoid neuron:

A neuron with sigmoid activation σ

$\sigma(w^T x)$ is usually interpreted as

the probability that target $z = 1$

Given a supervised dataset $X, y \quad y_i \in \{0, 1\} \quad \forall i$
sampled from $P_{\text{data}}(x, y)$, we can try to
learn target y via predicting

$$y = 1 \text{ if } \sigma(w^T x) > .5$$

$$y = 0 \text{ if } \sigma(w^T x) < .5$$

etc... maximum likelihood -- (recreate logistic regression)

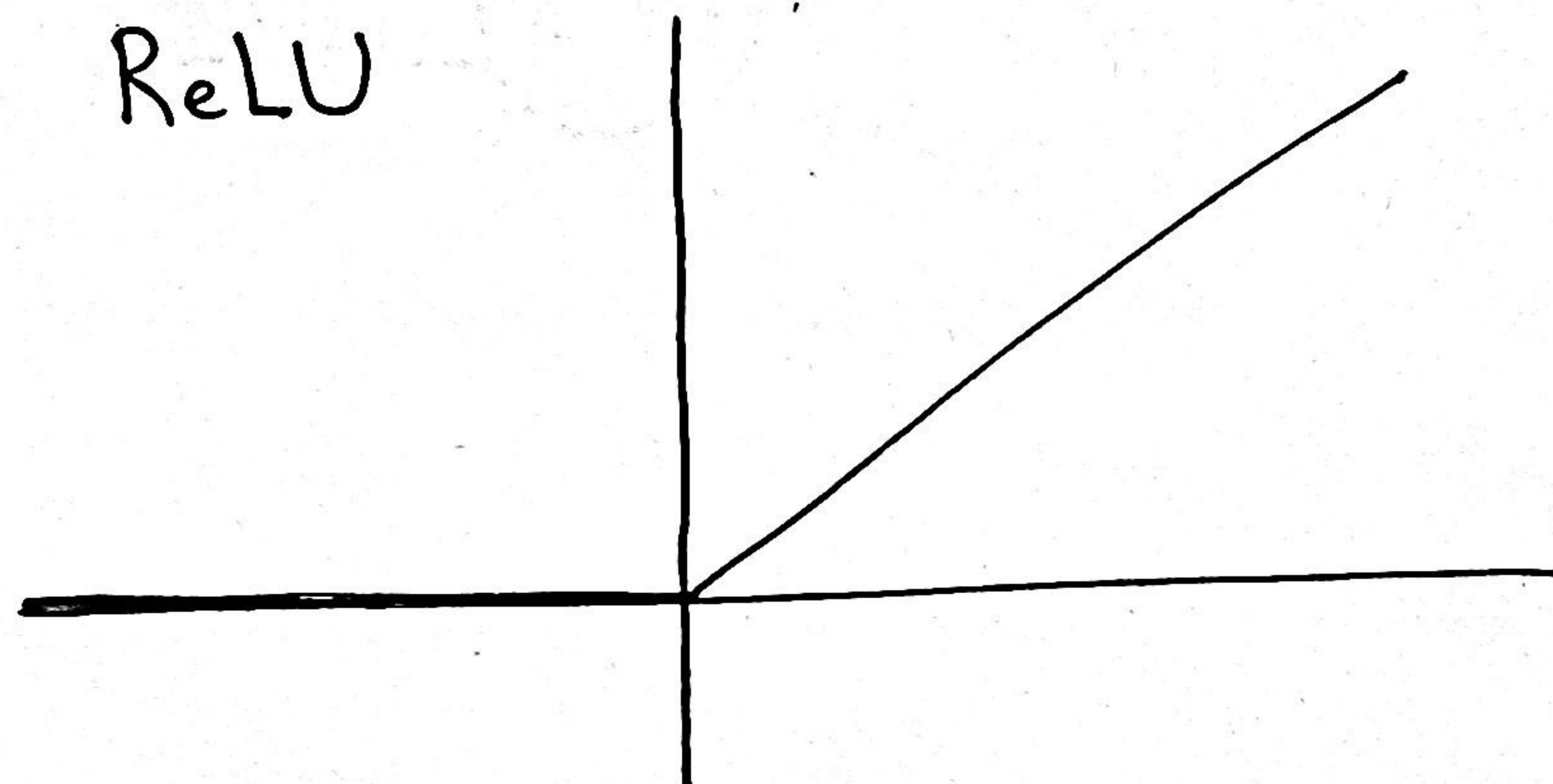
Linear Neurons: Neurons with activation $\phi(x) = x$

ReLU (rectified linear unit)

Activation function

$$\text{ReLU}(x) = x \quad H(x) = \begin{cases} 0 & x < 0 \\ x & x \geq 0 \end{cases}$$

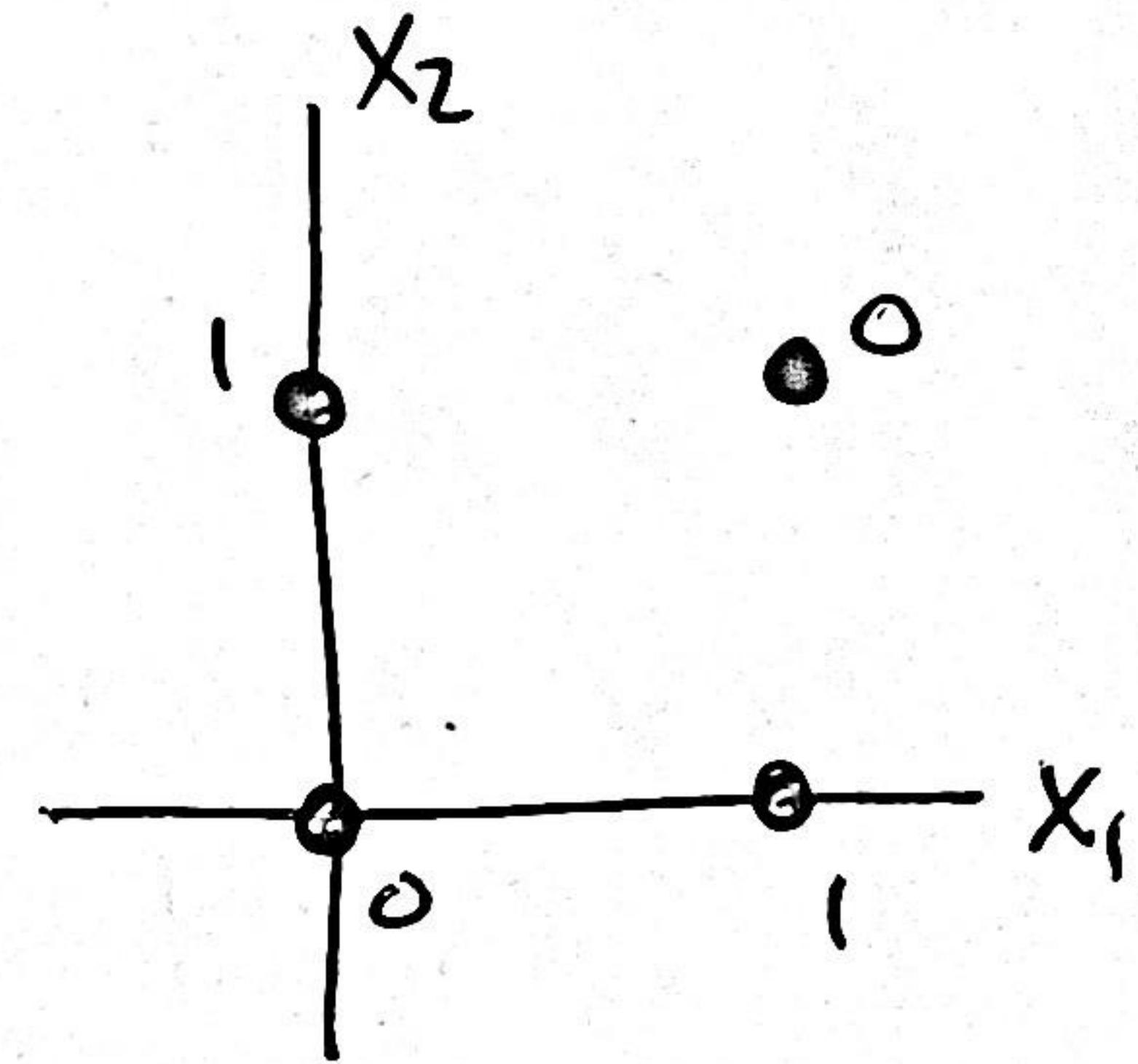
ReLU



A Network of Neurons:

E.g. $f^*(\vec{x}) = x_1 \text{ XOR } x_2$ $\vec{x} \in \{0,1\}^2$

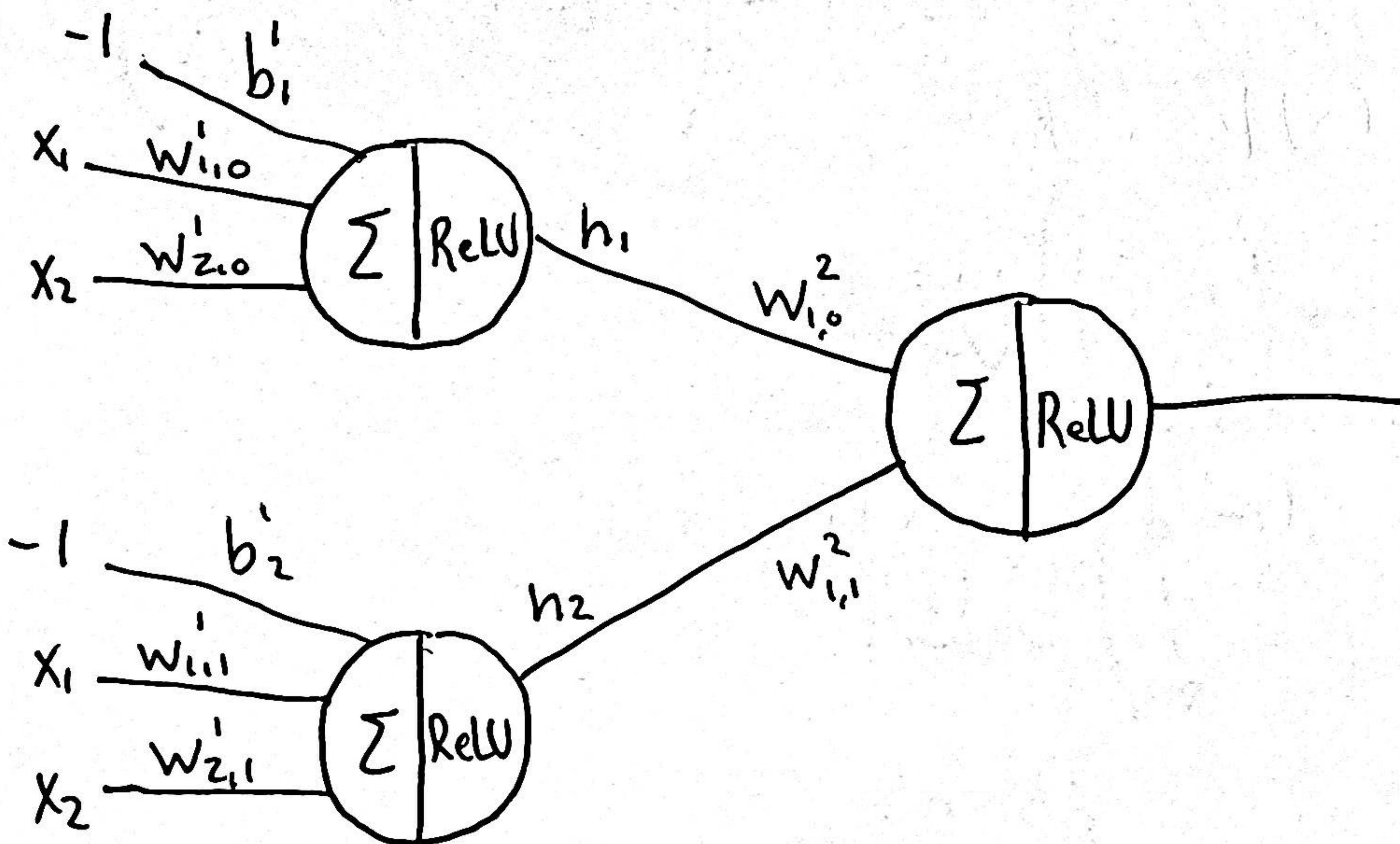
dataset $X = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}$ $y = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$



Consider a model $f(\vec{x}; \theta)$ and cost function

$$J(\theta) = \frac{1}{4} \sum_{\vec{x} \in X} (f^*(\vec{x}) - f(\vec{x}; \theta))^2$$

Idea



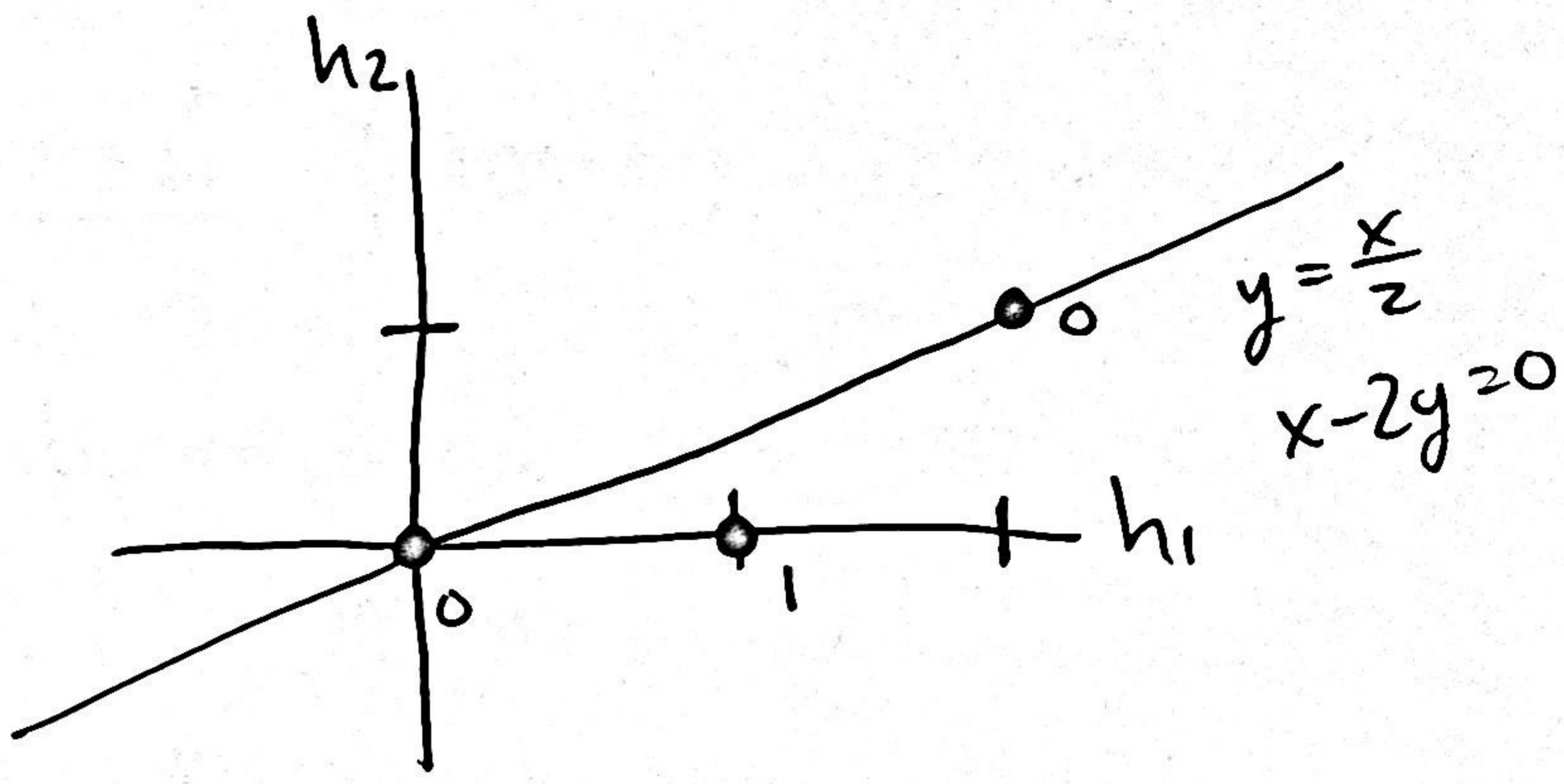
Idea: In x_1, x_2 -plane the points are not separable by a line but maybe by transforming $x_1, x_2 \rightarrow h_1, h_2$ the points can be separated by a line in the h_1, h_2 plane.

$$\begin{pmatrix} h_1 \\ h_2 \end{pmatrix} = \text{ReLU}\left(W^{(1)} \vec{x} + \vec{b}^{(1)}\right)$$

Try $W^{(1)} = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$ $\vec{b}^{(1)} = \begin{pmatrix} 0 \\ -1 \end{pmatrix}$

Apply to whole dataset

$$\begin{aligned} \begin{pmatrix} h_1 \\ h_2 \end{pmatrix} &= \text{ReLU}\left(XW^{(1)\top} + \vec{b}^{(1)\top}\right) = \text{ReLU}\left(\begin{pmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} + \begin{pmatrix} 0 & -1 \\ 0 & -1 \\ 0 & -1 \end{pmatrix}\right) \\ &= \text{ReLU}\left(\begin{pmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{pmatrix} + \begin{pmatrix} 0 & -1 \\ 0 & -1 \\ 0 & -1 \\ 0 & -1 \end{pmatrix}\right) \\ &= \text{ReLU}\left(\begin{pmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{pmatrix}\right) \\ &= \begin{pmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{pmatrix} \end{aligned}$$

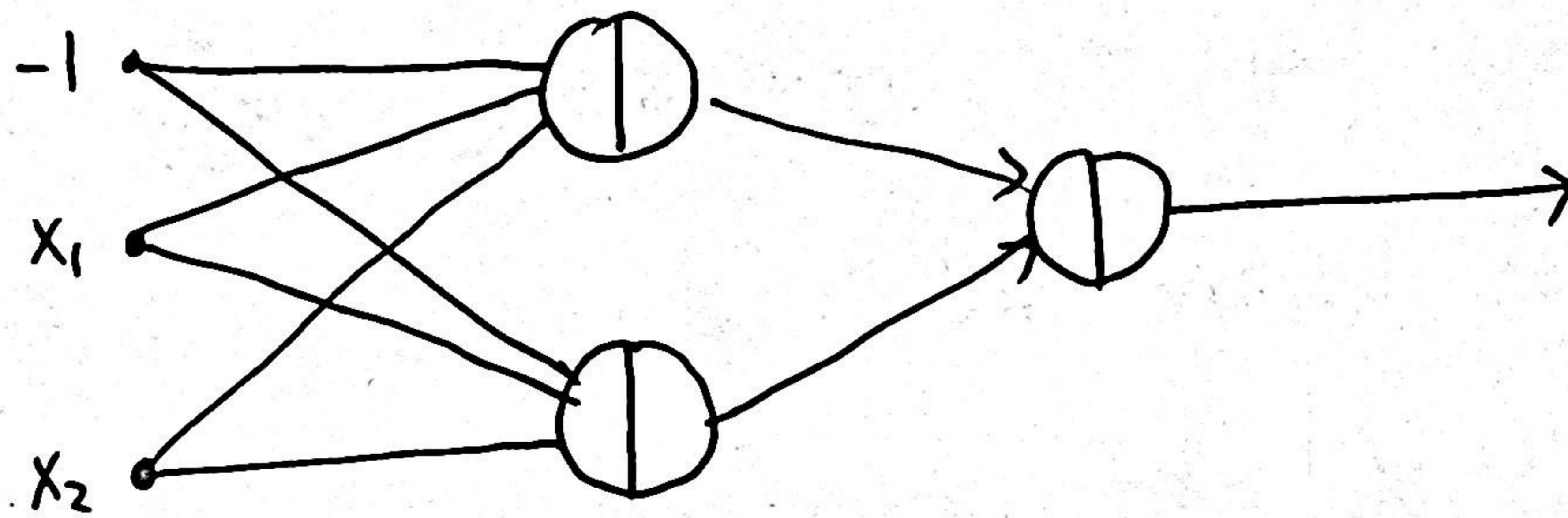


Final output is $W^{(2)} \begin{pmatrix} h_1 \\ h_2 \end{pmatrix} + b^{(2)}$

$$\text{Let } W^{(2)} = \begin{pmatrix} 1 \\ -2 \end{pmatrix} \quad b^{(2)} = 0$$

$$\text{ReLU} \left(\begin{pmatrix} 0 & 0 \\ 1 & 0 \\ 2 & 1 \end{pmatrix} \begin{pmatrix} 1 \\ -2 \end{pmatrix} \right) = \text{ReLU} \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$$

Better way to draw this:



Def: A partially ordered set (poset) is a pair (P, \leq) where P is a set and \leq is a relation on P s.t.

$$\textcircled{1} \quad x \leq x \quad \forall x \in P$$

$$\textcircled{2} \quad \text{if } x \leq y \text{ and } y \leq x \text{ then } x = y \quad \forall x, y \in P$$

$$\textcircled{3} \quad \text{if } x \leq y \text{ and } y \leq z \text{ then } x \leq z \quad \forall x, y, z \in P$$

Notation write $x < y$ if $x \leq y$ and $x \neq y$

Call (P, \leq) a total order if $\forall x, y \in P$

$x \leq y$ or $y \leq x$.

A cover relation is a pair (x, y) s.t.

$x \leq y$ and $\nexists z$ with $x < z < y$.

The Hasse diagram of a finite poset

(P, \leq) consists of a vertex for each $x \in P$ and a directed edge from x to y iff $x \leq y$

drawn from left to right

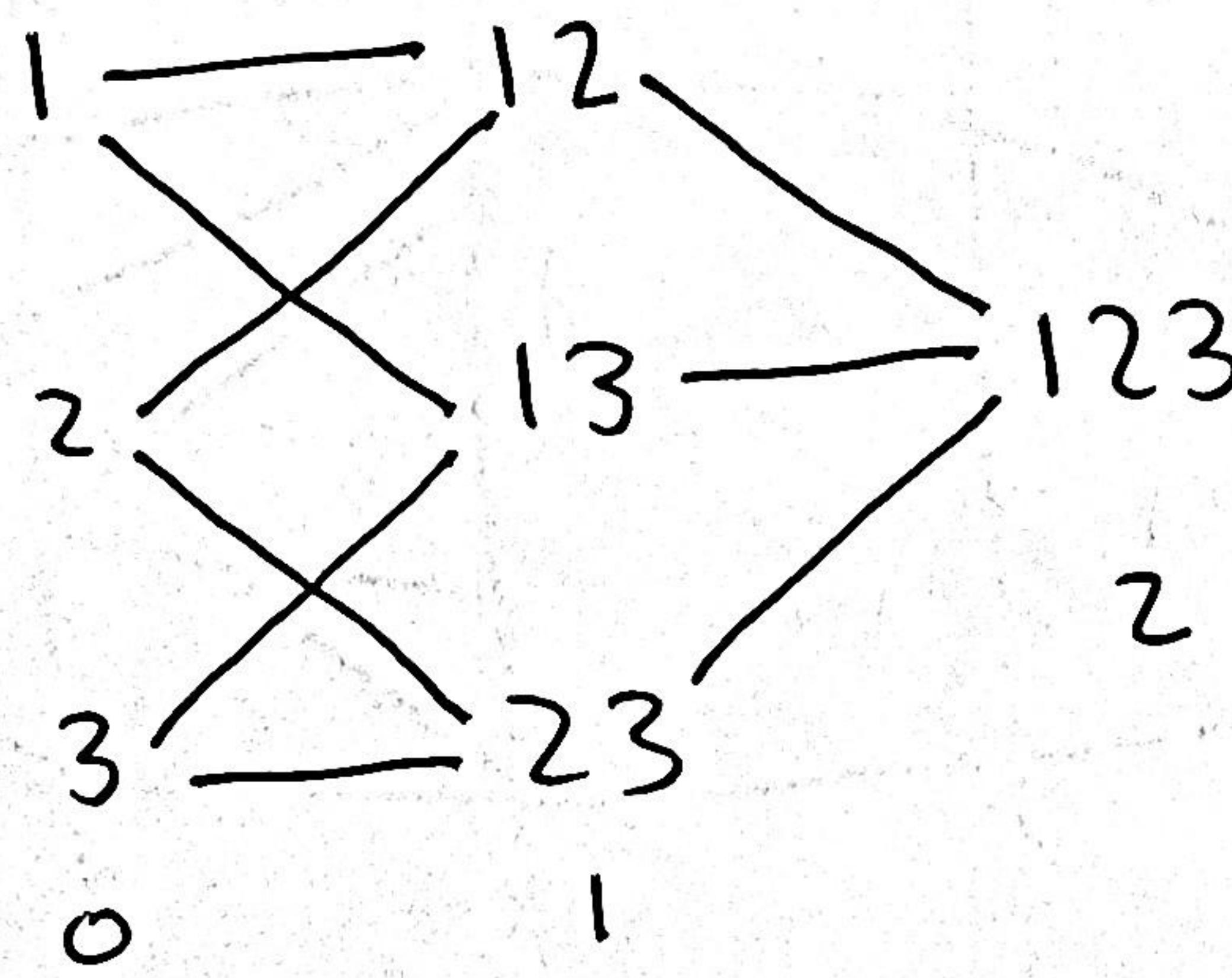
(so if $x \leq y$, draw x on the
left of y)

Write this as $H(P)$

E.g. $P = \{a, b, c, d\}$ $a \leq b \leq c \leq d$

$$\mathcal{H}(P) = \begin{array}{cccc} & \bullet & \rightarrow & \bullet & \rightarrow & \bullet & \rightarrow & \bullet \\ & a & & b & & c & & d \end{array}$$

E.g. $P = \{S \subseteq \{1, 2, 3\} \mid S \neq \emptyset\}$ ordered by containment

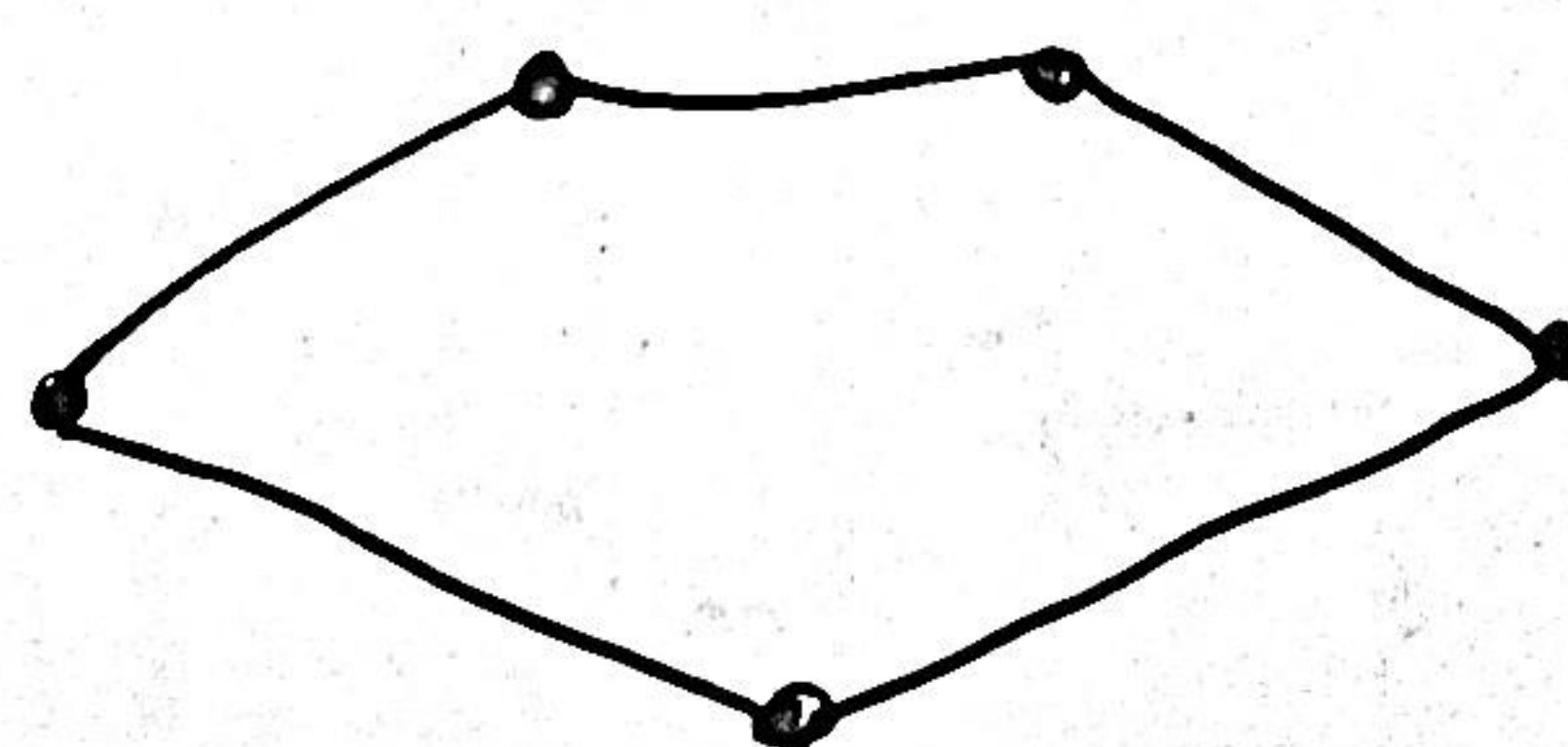


Def: A poset is graded if $\exists f: P \rightarrow \mathbb{N}$ (rank function)

- s.t.
- ① $x < y \Rightarrow f(x) < f(y) \quad \forall x, y \in P$
 - ② $x < y \Rightarrow f(x) + 1 = f(y) \quad \forall x, y \in P$

E.g. $\{S \subseteq \{1, 2, 3\} \mid S \neq \emptyset\}$ is graded

E.g.

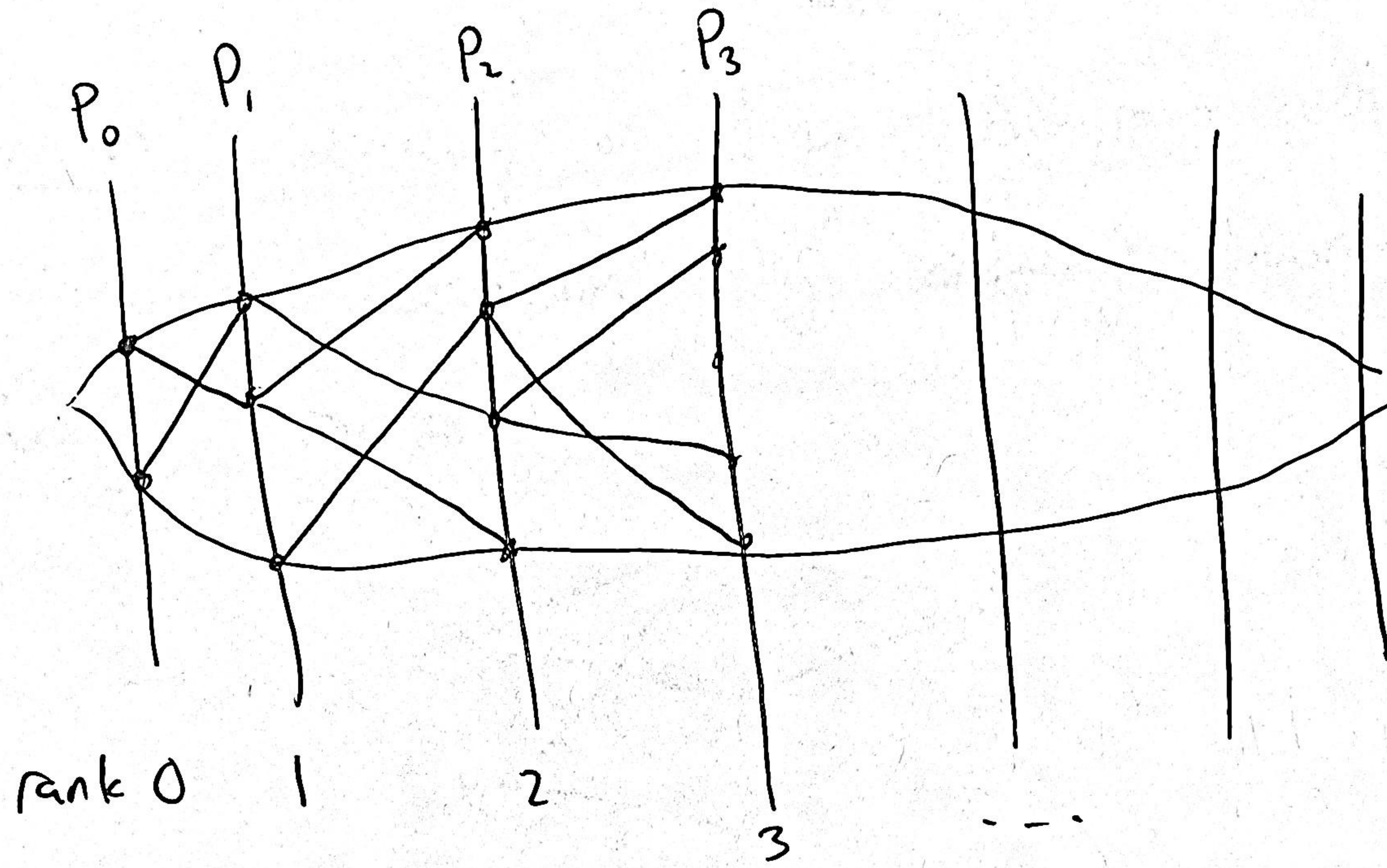


not graded

i.e. in a graded poset there is a notion of "depth"

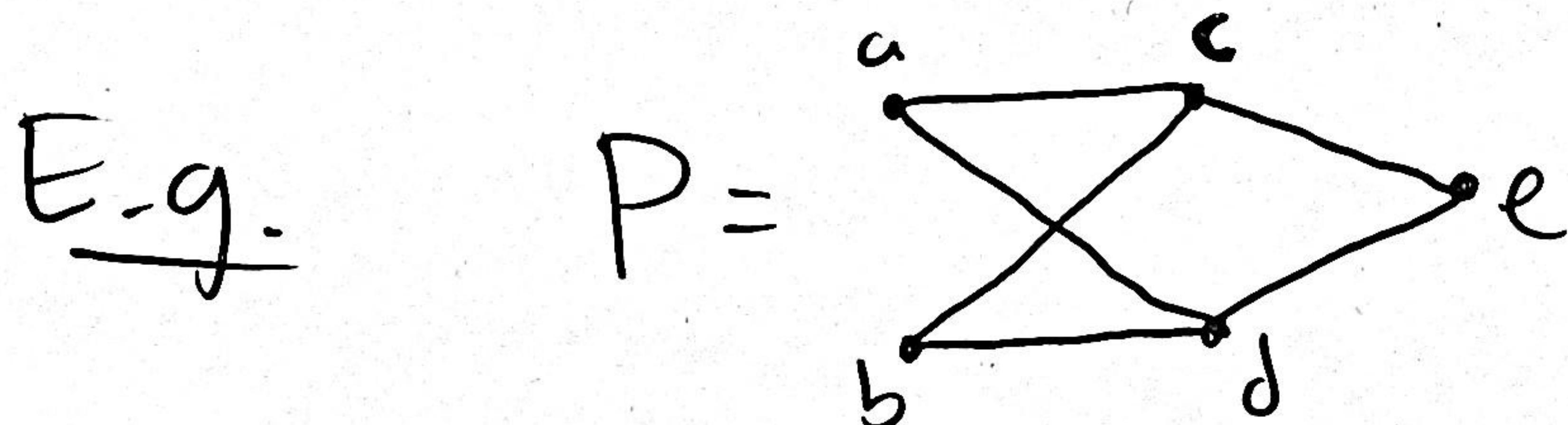
The rank decomposition of a graded poset ~~(P, ≤)~~
 (P, \leq, φ) is

$$P = \bigcup_{i \in N} P_i \quad P_i = \{x \in P \mid \varphi(x) = i\}$$



A linear extension of a poset (P, \leq) is
 a total ordering \leq' on P s.t.

$$x \leq y \implies x \leq' y \quad \forall x, y \in P$$

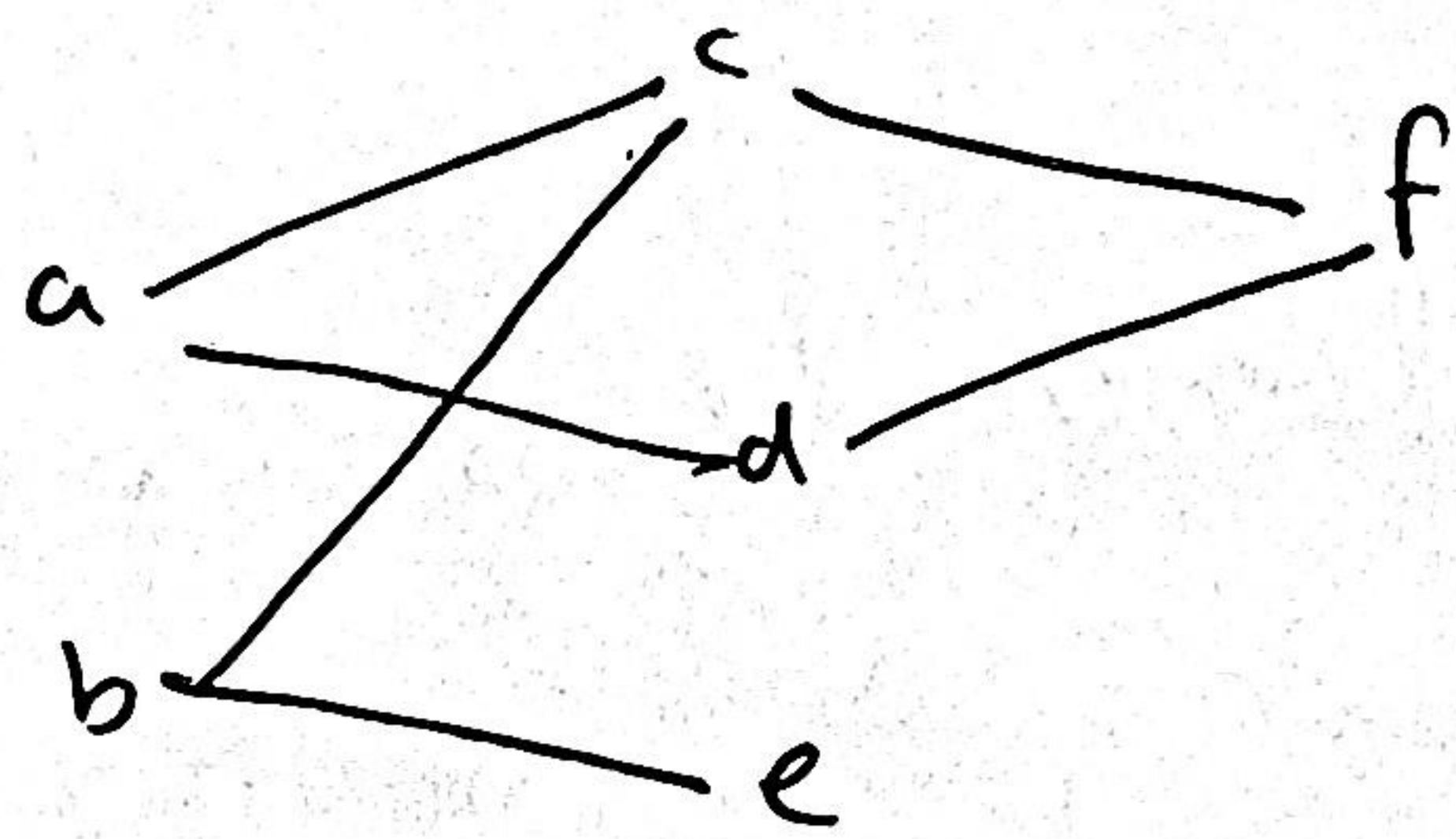


$$\leq': b < a < c < d < e$$

Lemma A linear extension of a graded poset is determined by a choice of total ordering of each of its ranks

\therefore There are $T \prod_{i \in N} |P_i|!$

Convention Drawing a Hasse diagram determines a linear extension by ordering downwards along ranks



$$a < b < c < d < e < f$$

\therefore Hasse diagram contains all information $(P, \leq, \dagger, \leq')$ (up to rank of P_0)

Def: A \mathbb{N} -graded vector space is a vector space V together with a decomposition

$$V = \bigoplus_{i \in \mathbb{N}} V_i \quad \left(\text{where } V_i \text{ is a subspace of } V \right)$$

E.g. $\mathbb{R}[x]$ is an \mathbb{N} -graded vector space

where $(\mathbb{R}[x])_i = \{cx^i \mid c \in \mathbb{R}\}$

Def: Given a graded poset (P, \leq, f, \leq') with linear extension \leq' , the associated graded vector space is

$$V_P^{\mathbb{N}} = \bigoplus_{i \in \mathbb{N}} V_P^{(i)}$$

$$V_P^{(i)} = \langle x \in P \mid f(x) = i \rangle \quad \left(\begin{array}{l} \text{real vector} \\ \text{space with} \\ \text{basis } P_i \end{array} \right)$$

basis P_i is an ordered basis with ordering specified by \leq'

Def: Given vector spaces V, V' , an affine map

$f: V \rightarrow V'$ is a map s.t.

$\forall \{x_1, \dots, x_n\} \subseteq V$ and $\forall \lambda_1, \dots, \lambda_n \in \mathbb{R}$ with $\sum_i \lambda_i = 1$

we have $f(\sum_i \lambda_i x_i) = \sum_i \lambda_i f(x_i)$

Lemma $f: V \rightarrow V'$ is affine iff \exists a linear

map $L: V \rightarrow V'$ and $\vec{b}_0 \in V'$ s.t.

$$f(v) = L(v) + \vec{b}_0 \quad \forall v \in V$$

Proof: Exercise

Note: After picking bases for V, V' an

affine map looks like

$$f(v) = Mv + \vec{b} \quad M \in \mathbb{R}^{\dim V' \times \dim V} \\ \vec{b} \in \mathbb{R}^{\dim V'}$$

Def Let $H = (P, \leq, f, \leq')$ be a finite graded poset with linear extension \leq' , s.t. $\min\{f(x) \mid x \in P\} = 0$.

A feedforward Neural Network ~~N~~ with architecture H consists of

- ① a collection $F = (F^{(0)}, \dots, F^{(D-1)})$ of affine maps where $D = \max\{f(x) \mid x \in P\}$
 s.t. $F^{(i)} : V_p^{(i)} \rightarrow V_p^{(i+1)}$ for $0 \leq i \leq D-1$
 and for $x \in V_p^{(i)}, y \in V_p^{(i+1)}$
 $\text{proj}_y F^{(i)}(x) = 0$ unless $x \leq y$

- ② a collection of maps

$$G = (G_1^{(1)}, \dots, G_D^{(D)}) \quad \text{where}$$

$$G_i^{(i)} : V_p^{(i)} \rightarrow V_p^{(i)}$$

Remark It is almost always the case that

$G^{(i)}$ is diagonal so $G^{(i)}(v_1, \dots, v_k) = (g_1^{(i)}(v_1), \dots, g_k^{(i)}(v_k))$
for some functions $g_j^{(i)} : \mathbb{R} \rightarrow \mathbb{R}$

Also usually $g_1^{(i)} = g_2^{(i)} = \dots = g_k^{(i)}$

Terminology ① D is called the depth of N

② ~~V_p~~ $V_p^{(i)}$ is the i^{th} layer of N

③ $V_p^{(1)}, \dots, V_p^{(D-1)}$ are called hidden layers

④ $V_p^{(D)}$ is the output layer

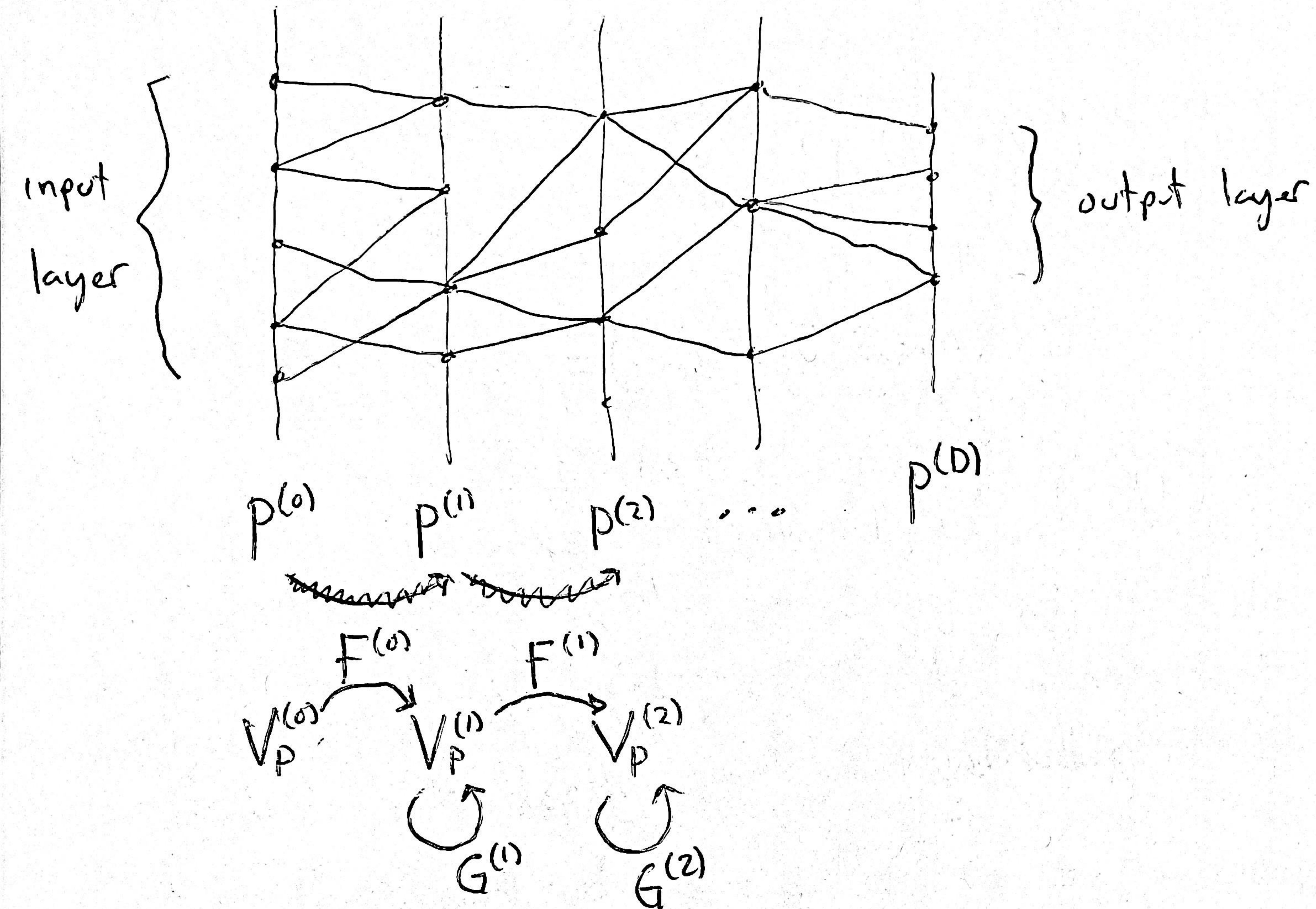
⑤ $V_p^{(0)}$ is the input layer

⑥ $G^{(i)}$ is the activation at the i^{th} layer

⑦ $g_j^{(i)}$ is the activation at the j^{th} neuron
of the i^{th} layer

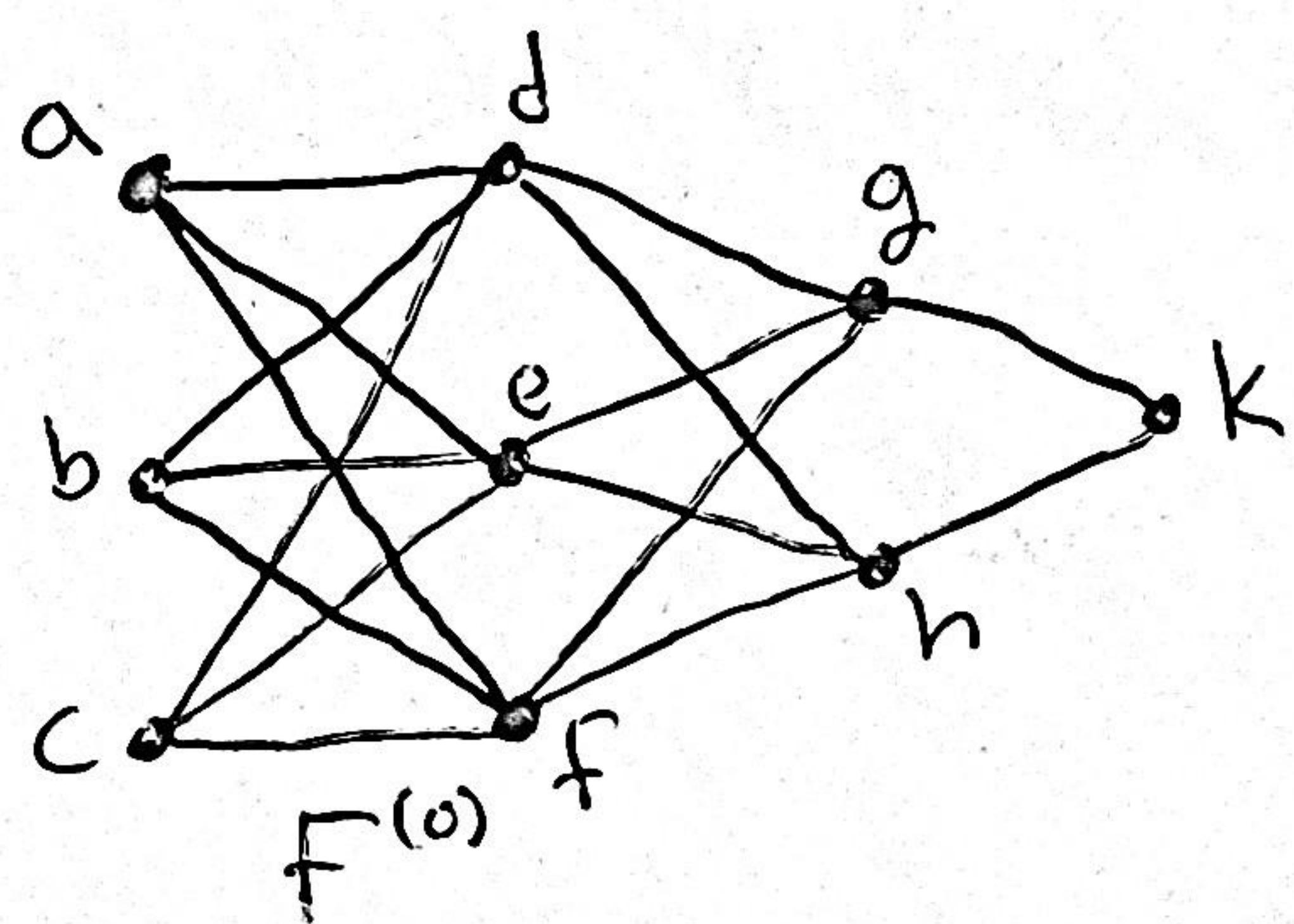
⑧ The matrix $M^{(i)}$ is the matrix of weights at layer i
 $\vec{b}^{(i)}$ is the bias at layer i

$$(\vec{f}_x^{(i)} = M^{(i)} \vec{x} + \vec{b}^{(i)})$$



E.g. Consider the poset $P(n_0, \dots, n_D)$ which has n_i elements at rank i for $0 \leq i \leq D$ and $x \leq y$ iff $f(y) = f(x) + 1$

$P(3, 3, 2, 1)$:



$$F^{(0)} \vec{x} = M^{(0)} \vec{x} + \vec{b}^{(0)} \quad M^{(0)} \in \mathbb{R}^{3 \times 3} \quad \vec{b}^{(0)} \in \mathbb{R}^3$$

$$F^{(1)} \vec{x} = M^{(1)} \vec{x} + \vec{b}^{(1)} \quad M^{(1)} \in \mathbb{R}^{2 \times 3} \quad \vec{b}^{(1)} \in \mathbb{R}^2$$

$$F^{(2)} \vec{x} = M^{(2)} \vec{x} + \vec{b}^{(2)} \quad M^{(2)} \in \mathbb{R}^{1 \times 2} \quad \vec{b}^{(2)} \in \mathbb{R}$$

~~Def~~

$$G^{(1)} \vec{x} = (g_{\varphi}^{(1)}(x_1), g_{\varphi}^{(1)}(x_2), g_{\varphi}^{(1)}(x_3))$$

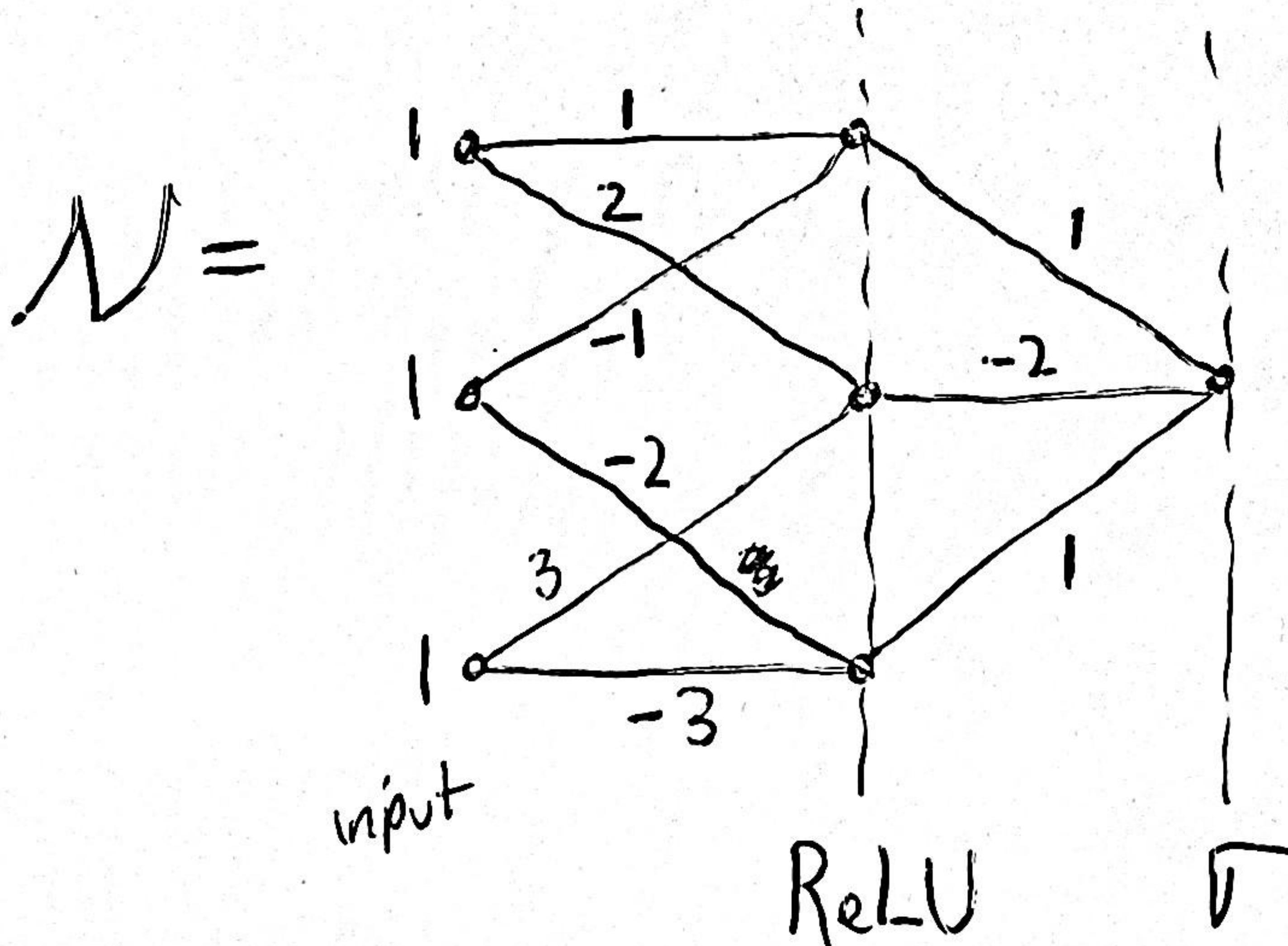
$$G^{(2)} \vec{x} = (g_{\varphi}^{(2)}(x_1), g_{\varphi}^{(2)}(x_2))$$

$$G^{(3)} \vec{x} = (g^{(3)}(x))$$

Def: Let $\mathcal{N} = (H, F, G)$ be a neural network with depth D and architecture $H = (P, \leq, \varphi, \leq')$. The feedforward of $\vec{x} \in \mathbb{R}^{\dim V_P^{(0)}}$ is

$$\mathcal{N} \vec{x} = G_D F_{D-1} \dots G_2 F_1 G_1 F_0 \vec{x} \in \mathbb{R}^{\dim V_P^{(D)}}$$

E.g.



$$\vec{b}^{(0)} = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \quad \vec{b}^{(1)} = 0$$

$$N\begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} = -G^{(0)}F^{(0)}G^{(0)}F^{(0)}\begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

$$F^{(0)}\begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 1 & -1 & 0 \\ 2 & 0 & 3 \\ 0 & -2 & -3 \end{pmatrix}\begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} + \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$
$$\vec{b}^{(0)}$$