

# Neural Network

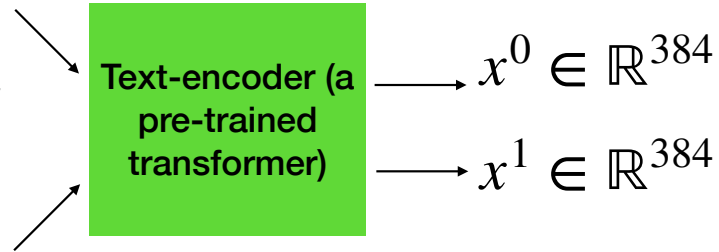
# Announcements

Kaggle competition will be released today

Task: sentiment analysis

r0: Bromwell High is a cartoon comedy. It ran at the same time as some other programs about school life, such as "Teachers". My 35 years in the teaching profession lead me to believe that Bromwell High's satire is much closer to reality than is "Teachers". The scramble to survive financially, the insightful students who can see right through their pathetic teachers' pomp, the pettiness of the whole situation, all remind me of the schools I knew and their students. When I saw the episode in which a student repeatedly tried to burn down the school, I immediately recalled ..... at ..... High. A classic line: INSPECTOR: I'm here to sack one of your teachers. STUDENT: Welcome to Bromwell High. I expect that many adults of my age think that Bromwell High is far fetched. What a pity that it isn't!

r1: Story of a man who has unnatural feelings for a pig. Starts out with a opening scene that is a terrific example of absurd comedy. A formal orchestra audience is turned into an insane, violent mob by the crazy chantings of it's singers. Unfortunately it stays absurd the WHOLE time with no general narrative eventually making it just too off putting. Even those from the era should be turned off. The cryptic dialogue would make Shakespeare seem easy to a third grader. On a technical level it's better than you might think with some good cinematography by future great Vilmos Zsigmond. Future stars Sally Kirkland and Frederic Forrest can be seen briefly.



Preference  $y \in \{0,1\}$  indicate which review is positive

# Recap on Boosting

Boosting iteratively learns a new classifier, and add it to the ensemble

Initialize  $H_1 = h_1 \in \mathcal{H}$

For  $t = 1 \dots$

# Recap on Boosting

Boosting iteratively learns a new classifier, and add it to the ensemble

Initialize  $H_1 = h_1 \in \mathcal{H}$

$$H_t = \sum_{i=1}^t \alpha_i h_i$$

For  $t = 1 \dots$

Denote  $\hat{\mathbf{y}} = [\underline{H_t(x_1)}, \underline{H_t(x_2)}, \dots, \underline{H_t(x_n)}]^\top \in \mathbb{R}^n$

# Recap on Boosting

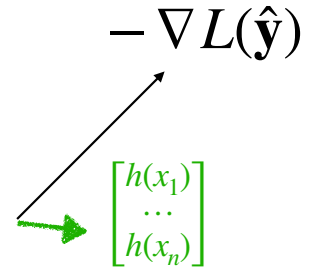
$$\nabla L(\hat{y}) = \begin{bmatrix} \frac{\partial l}{\partial \hat{y}_1} \\ \vdots \\ \frac{\partial l}{\partial \hat{y}_n} \end{bmatrix}$$

Boosting iteratively learns a new classifier, and add it to the ensemble

Initialize  $H_1 = h_1 \in \mathcal{H}$

For  $t = 1 \dots$

Denote  $\hat{y} = [H_t(x_1), H_t(x_2), \dots, H_t(x_n)]^T \in \mathbb{R}^n$



# Recap on Boosting

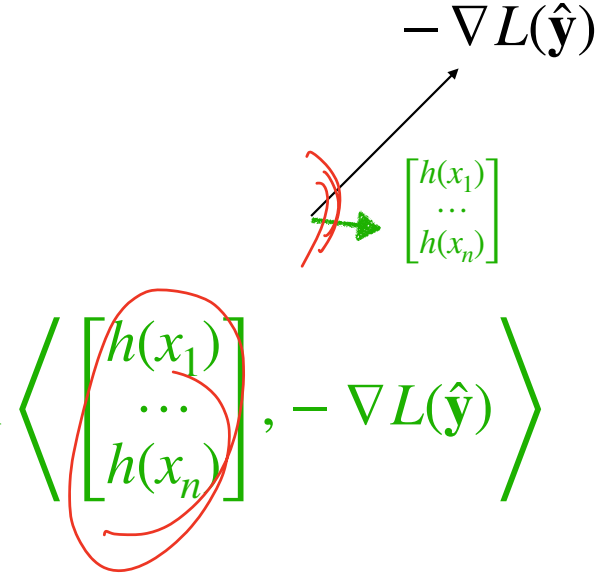
Boosting iteratively learns a new classifier, and add it to the ensemble

Initialize  $H_1 = h_1 \in \mathcal{H}$

For  $t = 1 \dots$

Denote  $\hat{\mathbf{y}} = [H_t(x_1), H_t(x_2), \dots, H_t(x_n)]^T \in \mathbb{R}^n$

Solve the optimization problem:  $h_{t+1} = \arg \max_{h \in \mathcal{H}}$


$$\left\langle \begin{bmatrix} h(x_1) \\ \dots \\ h(x_n) \end{bmatrix}, -\nabla L(\hat{\mathbf{y}}) \right\rangle$$

# Recap on Boosting

Boosting iteratively learns a new classifier, and add it to the ensemble

Initialize  $H_1 = h_1 \in \mathcal{H}$

For  $t = 1 \dots$

Denote  $\hat{\mathbf{y}} = [H_t(x_1), H_t(x_2), \dots, H_t(x_n)]^\top \in \mathbb{R}^n$

Solve the optimization problem:  $h_{t+1} = \arg \max_{h \in \mathcal{H}} \left\langle \begin{bmatrix} h(x_1) \\ \dots \\ h(x_n) \end{bmatrix}, -\nabla L(\hat{\mathbf{y}}) \right\rangle$

$H_{t+1} = H_t + \alpha h_{t+1}$

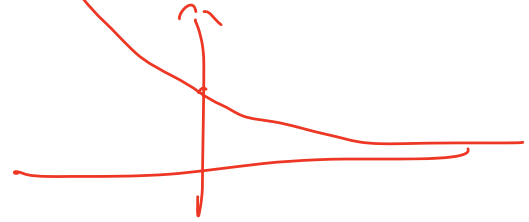
$-\nabla L(\hat{\mathbf{y}})$

$\begin{bmatrix} h(x_1) \\ \dots \\ h(x_n) \end{bmatrix}$

# Recap on AdaBoost

Adaboost follows this framework with  $\ell(\hat{y}, y) = \exp(-\hat{y} \cdot y)$

**1. Create a new weighted dataset:**





# Recap on AdaBoost

Adaboost follows this framework with  $\ell(\hat{y}, y) = \exp(-\hat{y} \cdot y)$

**1. Create a new weighted dataset:**

$$\hat{y}_i \leftarrow \text{He}(X_i)$$

For each  $x_i$ , compute  $p_i \propto \exp(-\hat{y}_i \cdot y_i)$

# Recap on AdaBoost

Adaboost follows this framework with  $\ell(\hat{y}, y) = \exp(-\hat{y} \cdot y)$

## 1. Create a new weighted dataset:

For each  $x_i$ , compute  $p_i \propto \exp(-\hat{y}_i \cdot y_i)$

Binary classification:  $h_{t+1} = \arg \min_{h \in \mathcal{H}} \sum_i p_i \cdot \mathbf{1}\{h(x_i) \neq y_i\}$

# Recap on AdaBoost

Adaboost follows this framework with  $\ell(\hat{y}, y) = \exp(-\hat{y} \cdot y)$

## 1. Create a new weighted dataset:

For each  $x_i$ , compute  $p_i \propto \exp(-\hat{y}_i \cdot y_i)$

Binary classification:  $h_{t+1} = \arg \min_{h \in \mathcal{H}} \sum_i p_i \cdot \mathbf{1}\{h(x_i) \neq y_i\}$

## 2. Add new learner to the ensemble:

# Recap on AdaBoost

Adaboost follows this framework with  $\ell(\hat{y}, y) = \exp(-\hat{y} \cdot y)$

## 1. Create a new weighted dataset:

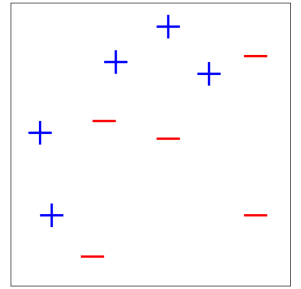
For each  $x_i$ , compute  $p_i \propto \exp(-\hat{y}_i \cdot y_i)$

Binary classification:  $h_{t+1} = \arg \min_{h \in \mathcal{H}} \sum_i p_i \cdot \mathbf{1}\{h(x_i) \neq y_i\}$

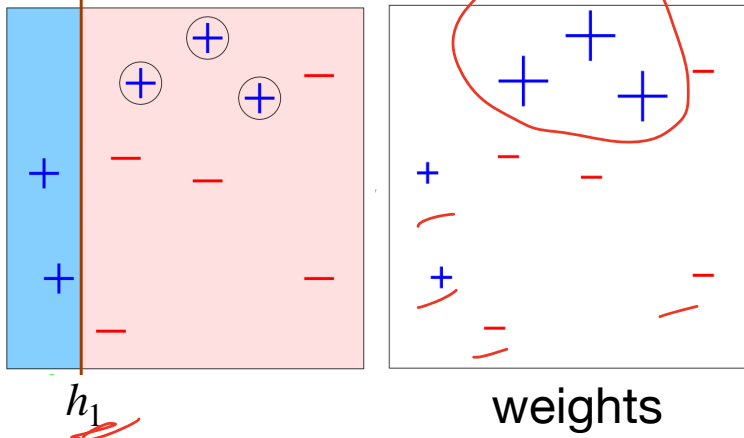
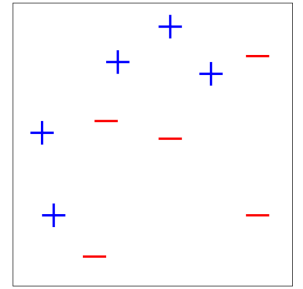
## 2. Add new learner to the ensemble:

$$H_{t+1} = H_t + \frac{1}{2} \ln \frac{1 - \epsilon}{\epsilon} \cdot h_{t+1}$$

Weaker learner: axis-aligned linear decision boundary

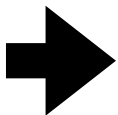
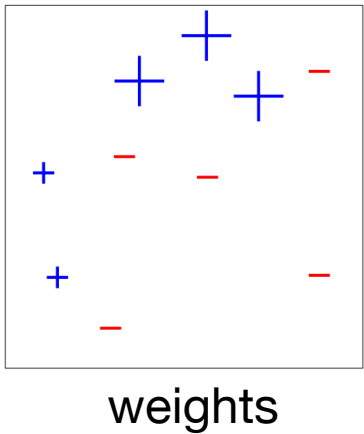
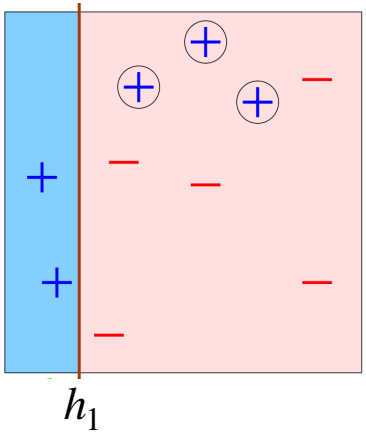


Weaker learner: axis-aligned linear decision boundary

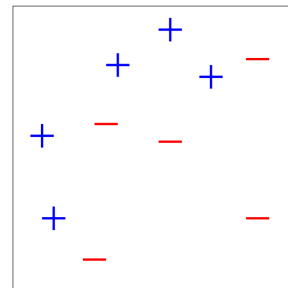


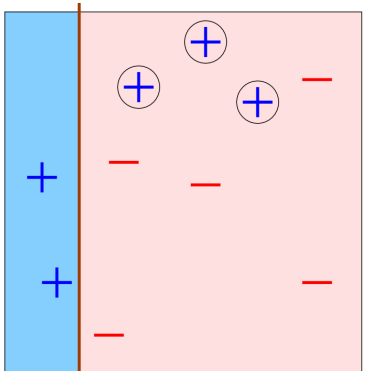
$$h_2 = h_1$$

$$P_i \propto \exp(-\hat{y}_i y_i)$$

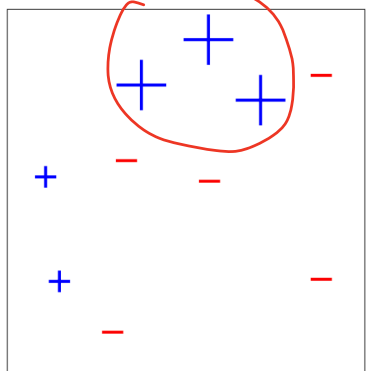


Weaker learner: axis-aligned linear decision boundary

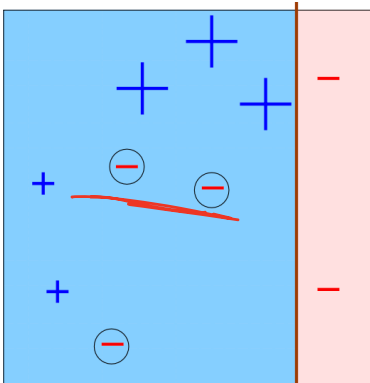
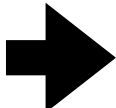




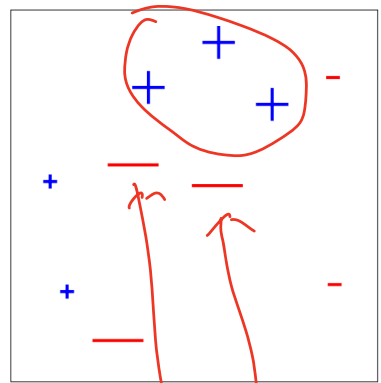
$h_1$



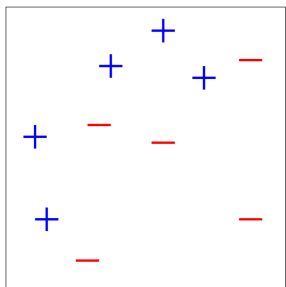
weights



$h_2$



Weaker learner: axis-aligned linear decision boundary

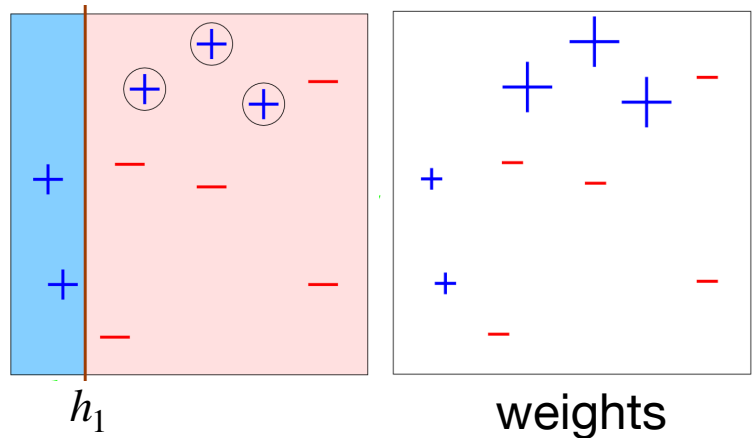


$$H_2 = h_1 + 2h_2$$

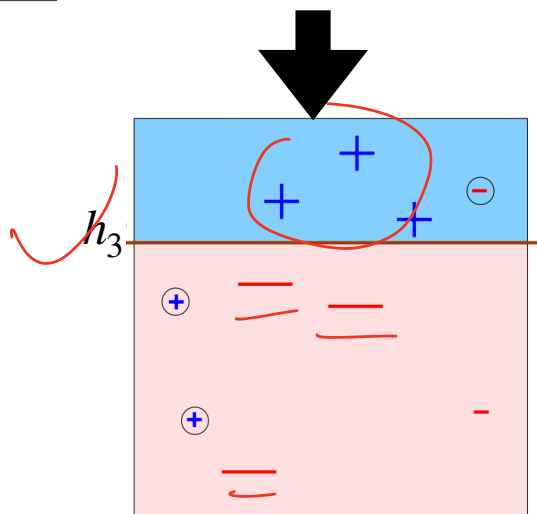
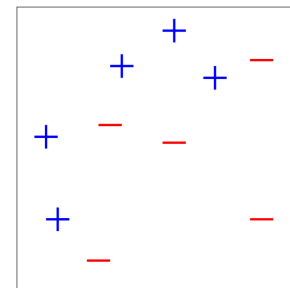
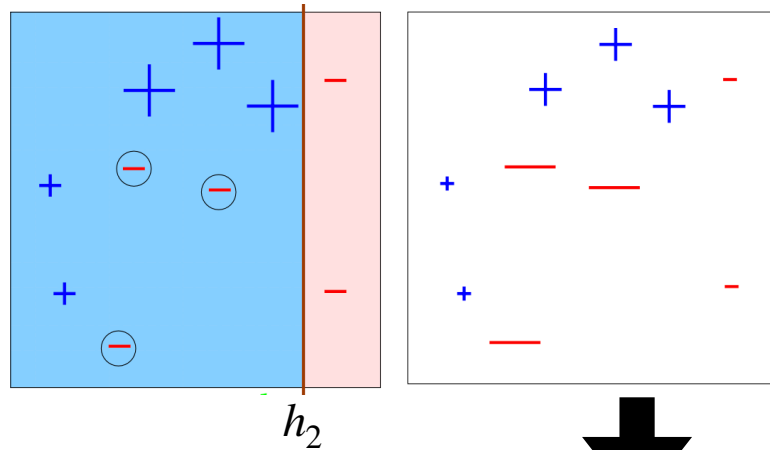
$$\hat{y}_i = H_2(x_i)$$

$$P_i \propto \exp(-\hat{y}_i y_i)$$

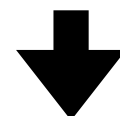
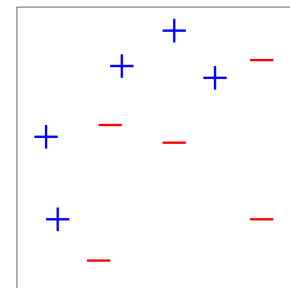
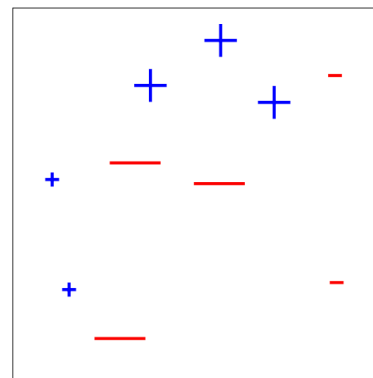
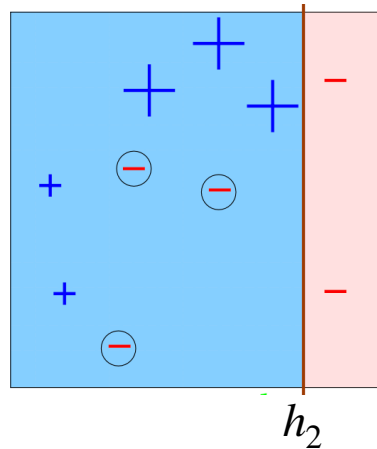
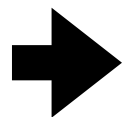
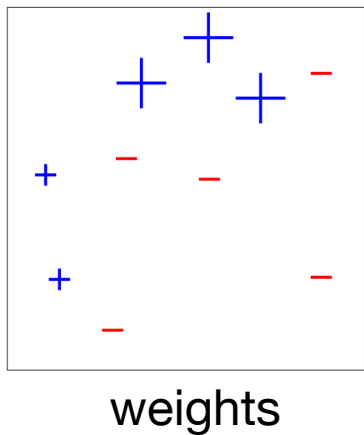
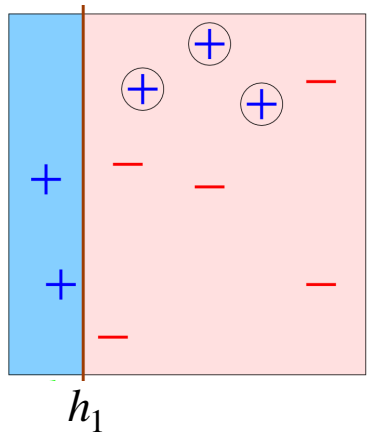




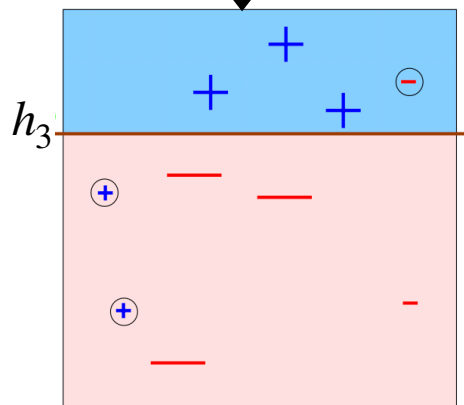
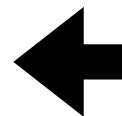
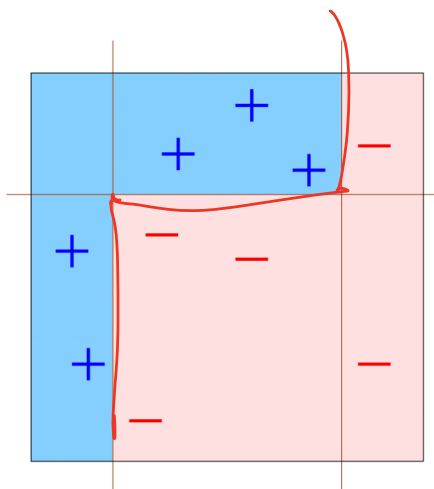
Weaker learner: axis-aligned linear decision boundary



Weaker learner: axis-aligned linear decision boundary



Final learner



# Outline of Today

1. Analysis of Boosting

2. Multilayer feedforward Neural Network

# The definition of Weak learning

Each weaker learning optimizes its own data:

$$\tilde{\mathcal{D}} = \{p_i, x_i, y_i\}, \text{ where } \sum_i p_i = 1, p_i \geq 0, \forall i$$

$$h_{t+1} = \arg \min_{h \in \mathcal{H}} \sum_{i=1}^n p_i \cdot \mathbf{1}(h(x_i) \neq y_i) \quad \checkmark$$

# The definition of Weak learning

Each weaker learning optimizes its own data:

$$\tilde{\mathcal{D}} = \{p_i, x_i, y_i\}, \text{ where } \sum_i p_i = 1, p_i \geq 0, \forall i$$

$$h_{t+1} = \arg \min_{h \in \mathcal{H}} \sum_{i=1}^n p_i \cdot \mathbf{1}(h(x_i) \neq y_i)$$

Assume that weaker learner's loss  $\epsilon := \sum_{i=1}^n p_i \mathbf{1}\{h_{t+1}(x_i) \neq y_i\} \leq \frac{1}{2} + \gamma, \gamma > 0$

# The definition of Weak learning

Each weaker learning optimizes its own data:

$$\tilde{\mathcal{D}} = \{p_i, x_i, y_i\}, \text{ where } \sum_i p_i = 1, p_i \geq 0, \forall i$$

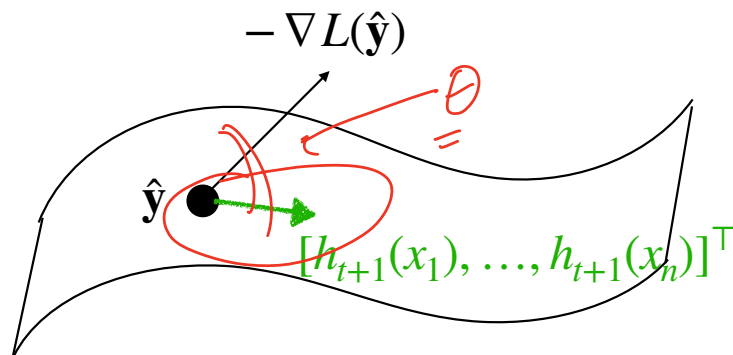
$$h_{t+1} = \arg \min_{h \in \mathcal{H}} \sum_{i=1}^n p_i \cdot \mathbf{1}(h(x_i) \neq y_i)$$

Assume that weaker learner's loss  $\epsilon := \sum_{i=1}^n p_i \mathbf{1}\{h_{t+1}(x_i) \neq y_i\} \leq \frac{1}{2} - \gamma, \gamma > 0$

Q: assume  $\mathcal{H}$  is symmetric. i.e.,  $h \in \mathcal{H}$  iff  $-h \in \mathcal{H}$ , why does the above always hold?

# Weaker learnability implies approximating gradient well

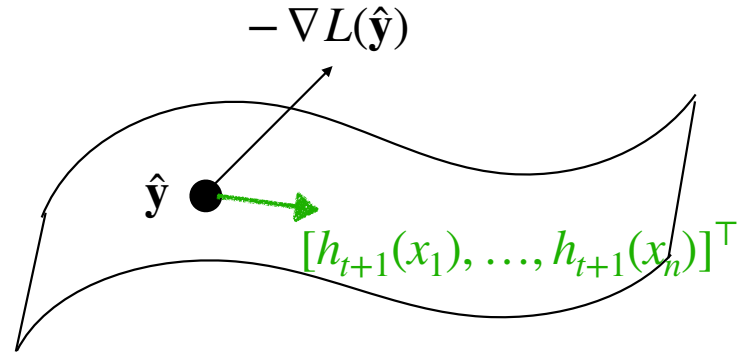
Assume that weaker learner's loss  $\epsilon := \sum_{i=1}^n p_i \mathbf{1}\{h_{t+1}(x_i) \neq y_i\} \leq \frac{1}{2} - \gamma$ ,  $\gamma > 0$



# Weaker learnability implies approximating gradient well

Assume that weaker learner's loss  $\epsilon := \sum_{i=1}^n p_i \mathbf{1}\{h_{t+1}(x_i) \neq y_i\} \leq \frac{1}{2} - \gamma$ ,  $\gamma > 0$

$$(-\nabla L(\hat{\mathbf{y}}))^\top \begin{bmatrix} h_{t+1}(x_1) \\ \dots \\ h_{t+1}(x_n) \end{bmatrix}$$

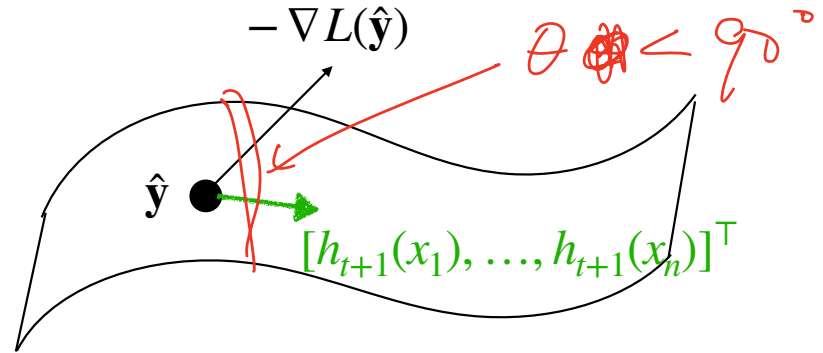




# Weaker learnability implies approximating gradient well

Assume that weaker learner's loss  $\epsilon := \sum_{i=1}^n p_i \mathbf{1}\{h_{t+1}(x_i) \neq y_i\} \leq \frac{1}{2} - \gamma$ ,  $\gamma > 0$

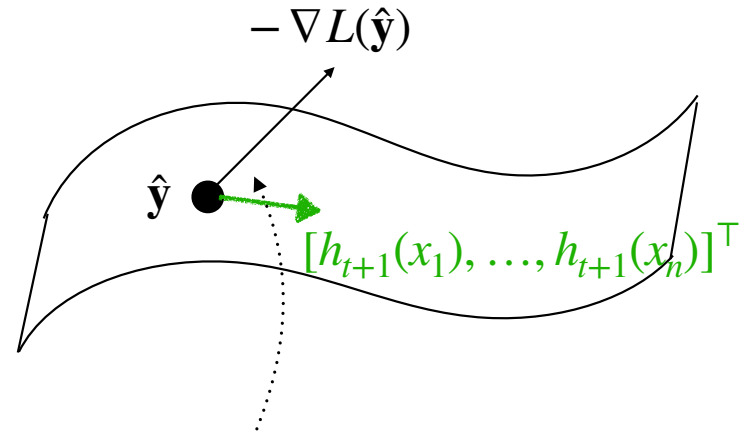
$$\begin{aligned} & (-\nabla L(\hat{y}))^\top \begin{bmatrix} h_{t+1}(x_1) \\ \dots \\ h_{t+1}(x_n) \end{bmatrix} \\ & \geq \left( \sum_{j=1}^n w_j \right) 2\gamma > 0 \end{aligned}$$



# Weaker learnability implies approximating gradient well

Assume that weaker learner's loss  $\epsilon := \sum_{i=1}^n p_i \mathbf{1}\{h_{t+1}(x_i) \neq y_i\} \leq \frac{1}{2} - \gamma$ ,  $\gamma > 0$

$$\begin{aligned} & (-\nabla L(\hat{y}))^\top \begin{bmatrix} h_{t+1}(x_1) \\ \dots \\ h_{t+1}(x_n) \end{bmatrix} \\ & \geq \left( \sum_{j=1}^n w_j \right) 2\gamma > 0 \end{aligned}$$



Within 90 degree, so  
improve the objective!

$$l(\hat{y}, y) = \exp(-\hat{y} \cdot y)$$

## Formal Convergence of AdaBoost

Then after  $T$  iterations, for the original exp loss, we have

$$\frac{1}{n} \sum_{i=1}^n \exp(-H_T(x_i) \cdot y_i) \leq n \underbrace{(1 - 4\gamma^2)}_{< 1}^{T/2}$$

$$\gamma \leftrightarrow \frac{1}{2} - \gamma$$

$$(0.75)^{\frac{100}{2}}$$

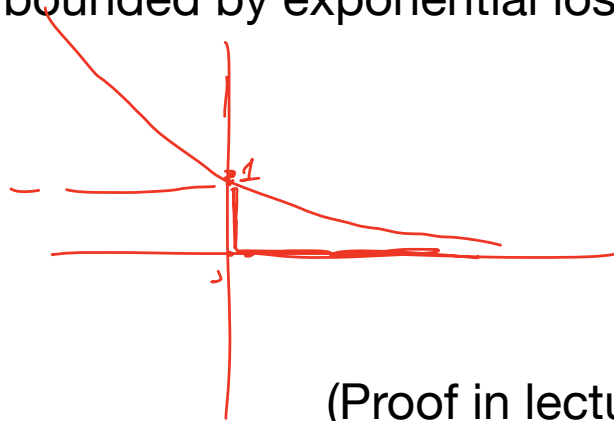
(Proof in lecture note, optional)

# Formal Convergence of AdaBoost

Then after  $T$  iterations, for the original exp loss, we have

$$\frac{1}{n} \sum_{i=1}^n \exp(-H_T(x_i) \cdot y_i) \leq n(1 - 4\gamma^2)^{T/2}$$

Note zero-one loss is upper bounded by exponential loss



(Proof in lecture note, optional)

# Formal Convergence of AdaBoost

Then after  $T$  iterations, for the original exp loss, we have

$$\frac{1}{n} \sum_{i=1}^n \exp(-H_T(x_i) \cdot y_i) \leq n(1 - 4\gamma^2)^{T/2}$$

Note zero-one loss is upper bounded by exponential loss

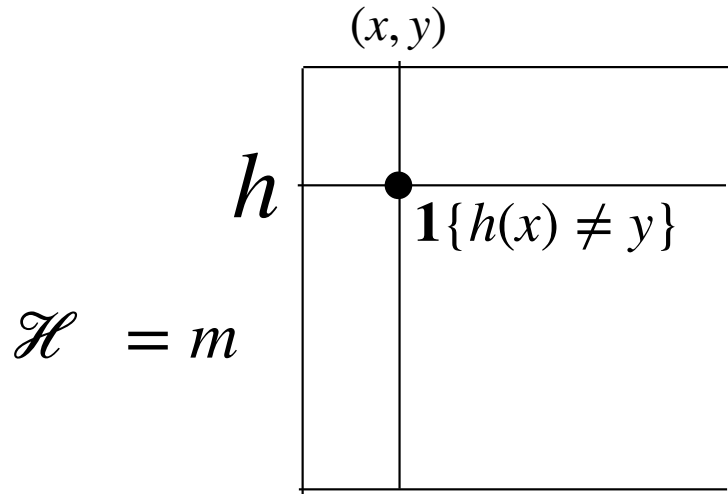
$$\frac{1}{n} \sum_{i=1}^n \mathbf{1}\{\text{sign}(H_T(x_i)) \neq y_i\} \leq \frac{1}{n} \sum_{i=1}^n \exp(-H_T(x_i) \cdot y_i) \leq n(1 - 4\gamma^2)^{T/2}$$

*Avg # of mistakes*

(Proof in lecture note, optional)

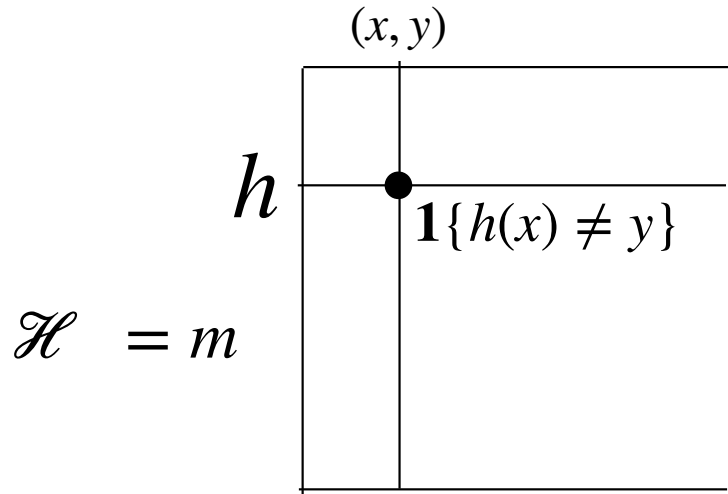
# Thinking about Boosting via two player zero sum game

$$\mathcal{D} = n$$



# Thinking about Boosting via two player zero sum game

$$\mathcal{D} = n$$

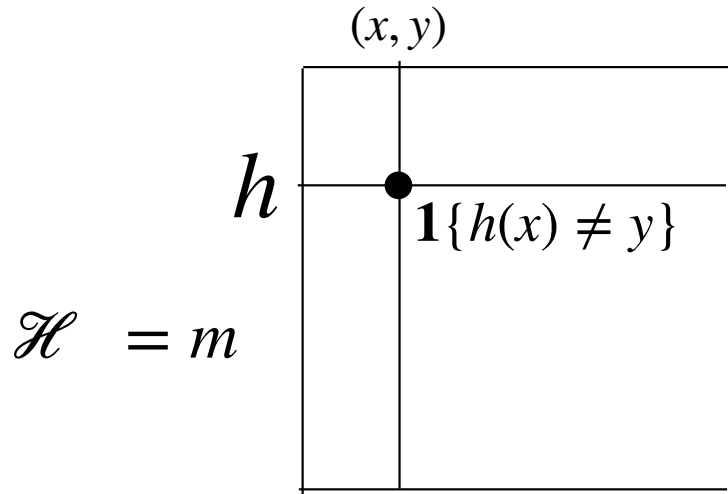


Row player plays hypothesis  $h \in \mathcal{H}$

Column player plays example  $(x, y)$

# Thinking about Boosting via two player zero sum game

$$\mathcal{D} = n$$



Row player plays hypothesis  $h \in \mathcal{H}$

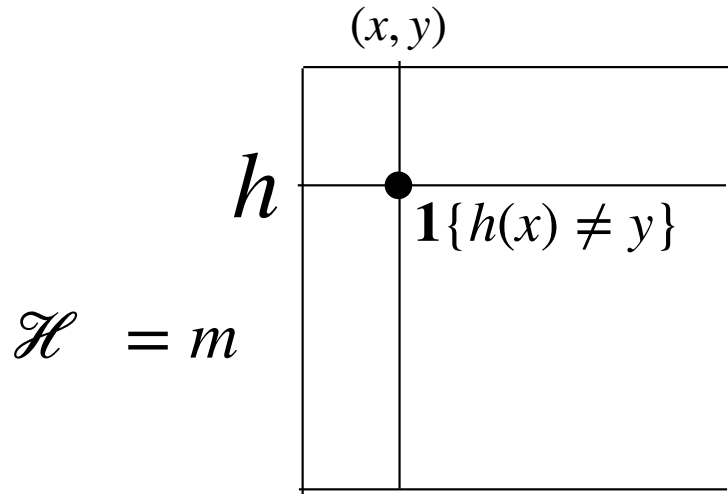
Column player plays example  $(x, y)$

Row player gets loss  $\mathbf{1}\{h(x) \neq y\}$



# Thinking about Boosting via two player zero sum game

$$\mathcal{D} = n$$



Row player plays hypothesis  $h \in \mathcal{H}$

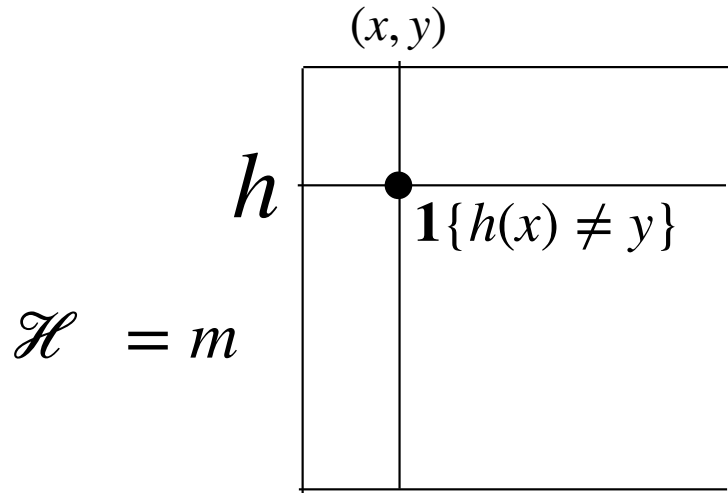
Column player plays example  $(x, y)$

Row player gets loss  $\mathbf{1}\{h(x) \neq y\}$

Column player gets loss  $-\mathbf{1}\{h(x) \neq y\}$

# Thinking about Boosting via two player zero sum game

$$\mathcal{D} = n$$



Row player plays hypothesis  $h \in \mathcal{H}$

Column player plays example  $(x, y)$

Row player gets loss  $\mathbf{1}\{h(x) \neq y\}$

Column player gets loss  $-\mathbf{1}\{h(x) \neq y\}$

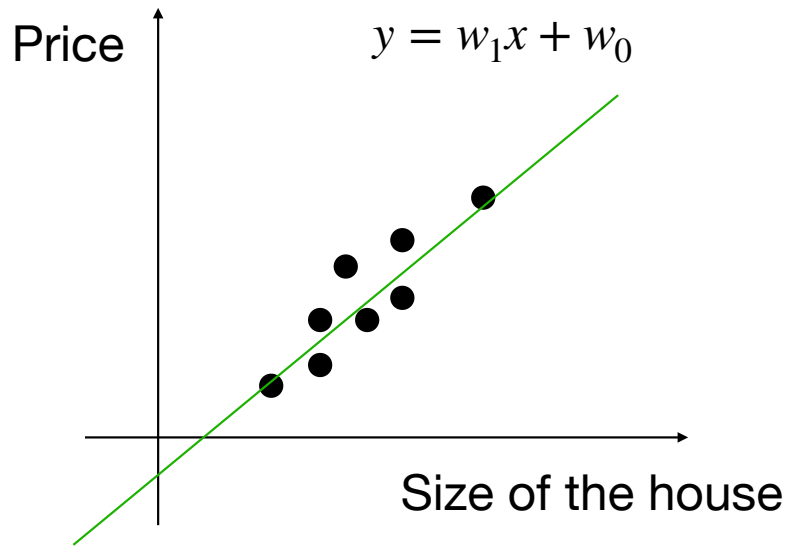
Boosting can be understood as running some specific algorithm to find the Nash equilibrium of the game

# Outline of Today

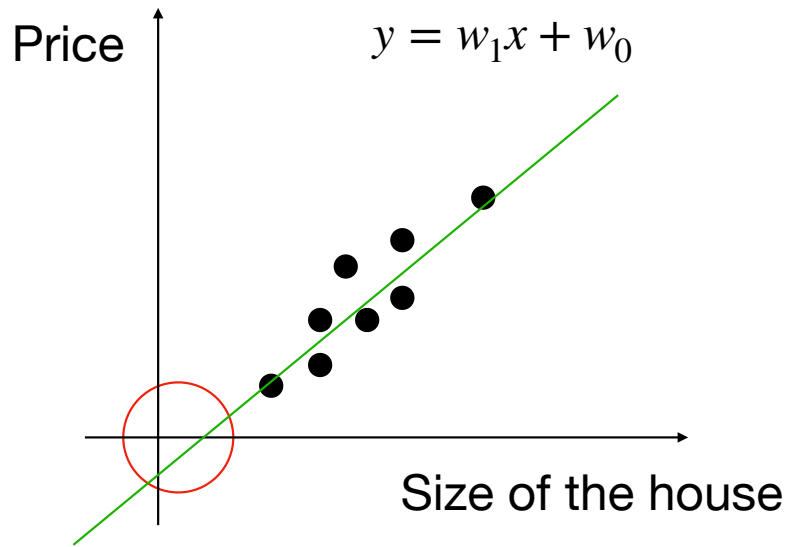
1. Analysis of Boosting

2. Multilayer feedforward Neural Network

# Linear Regression Revisit

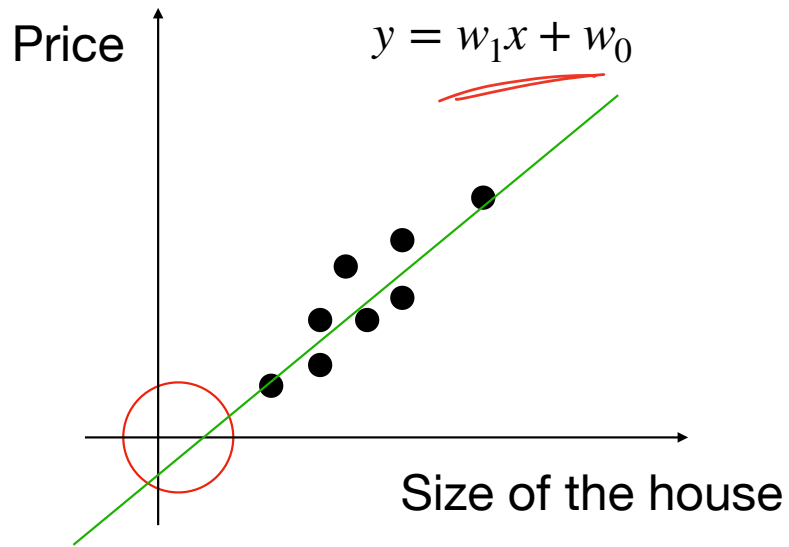


# Linear Regression Revisit



Negative part does not  
make too much sense

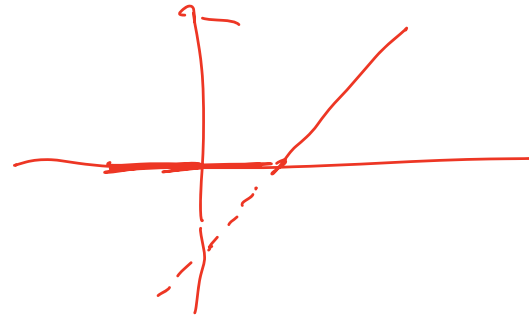
# Linear Regression Revisit



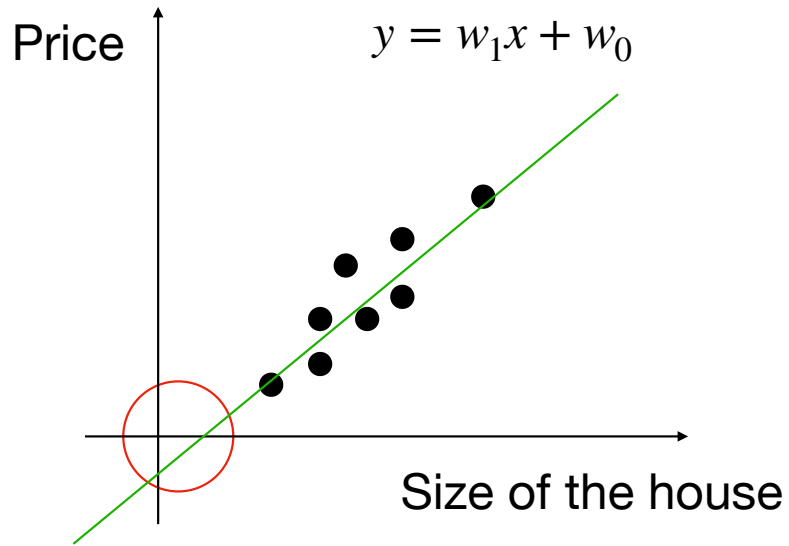
Negative part does not  
make too much sense

We can fix this with a simple  
nonlinear function

$$y = \max\{\underline{w_1x + w_0}, 0\}$$



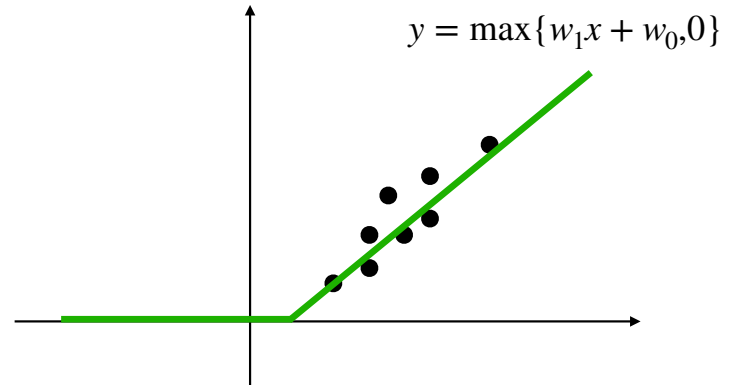
# Linear Regression Revisit



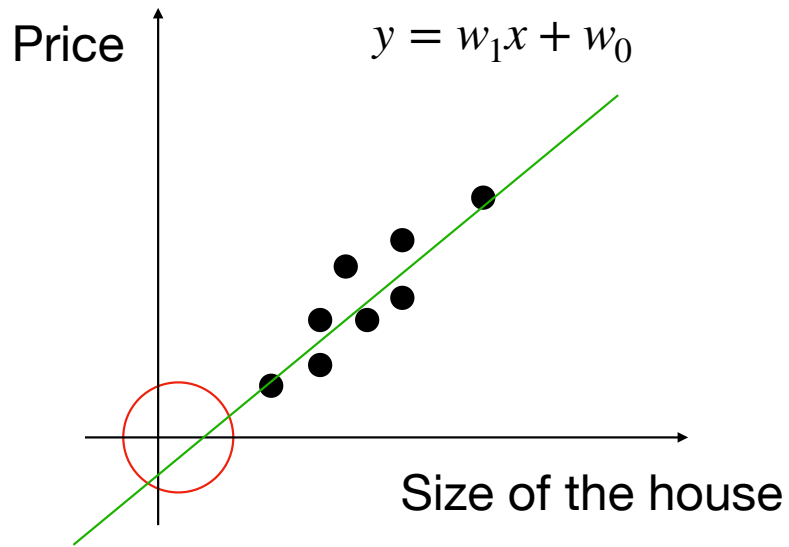
Negative part does not make too much sense

We can fix this with a simple nonlinear function

$$y = \max\{w_1x + w_0, 0\}$$



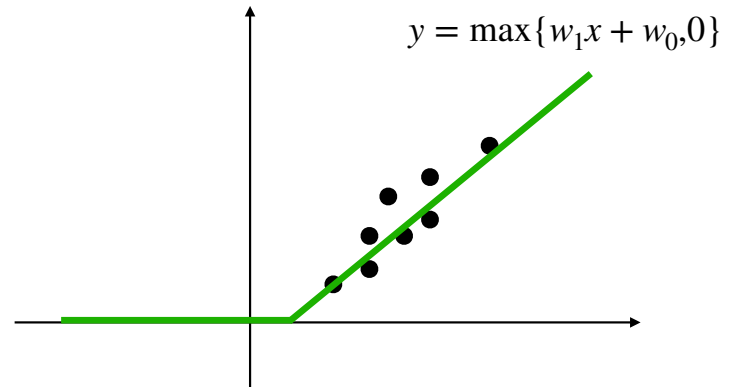
# Linear Regression Revisit



Negative part does not make too much sense

We can fix this with a simple nonlinear function

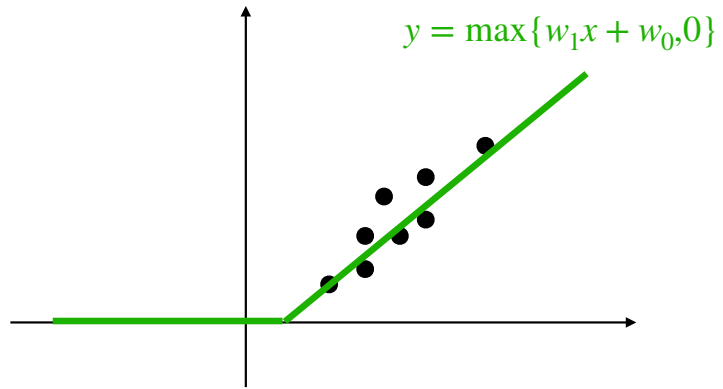
$$y = \max\{w_1x + w_0, 0\}$$



rectified linear unit (ReLU)

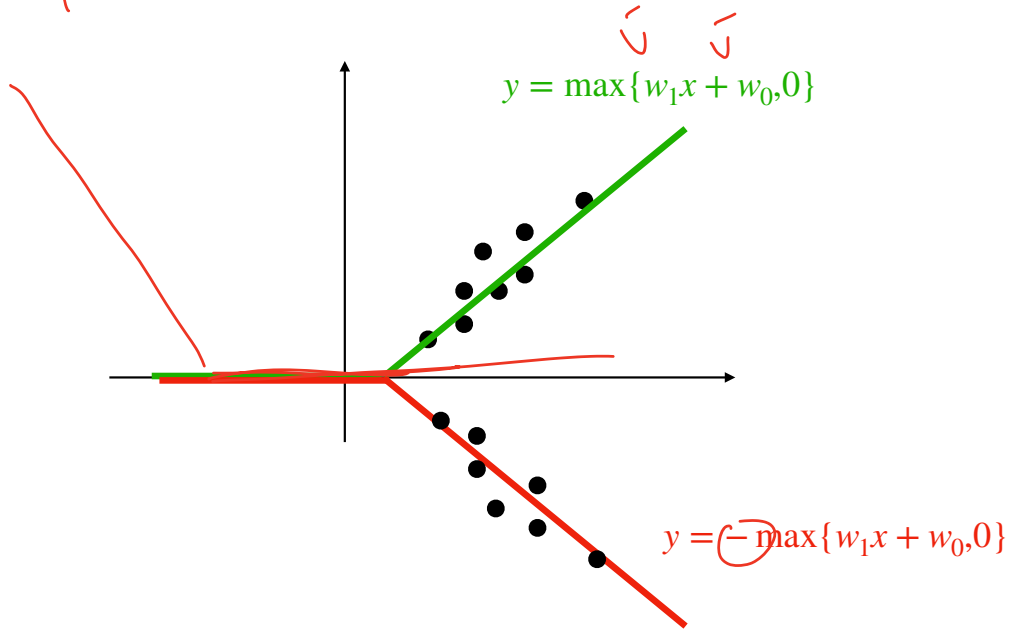


# A single neuron network

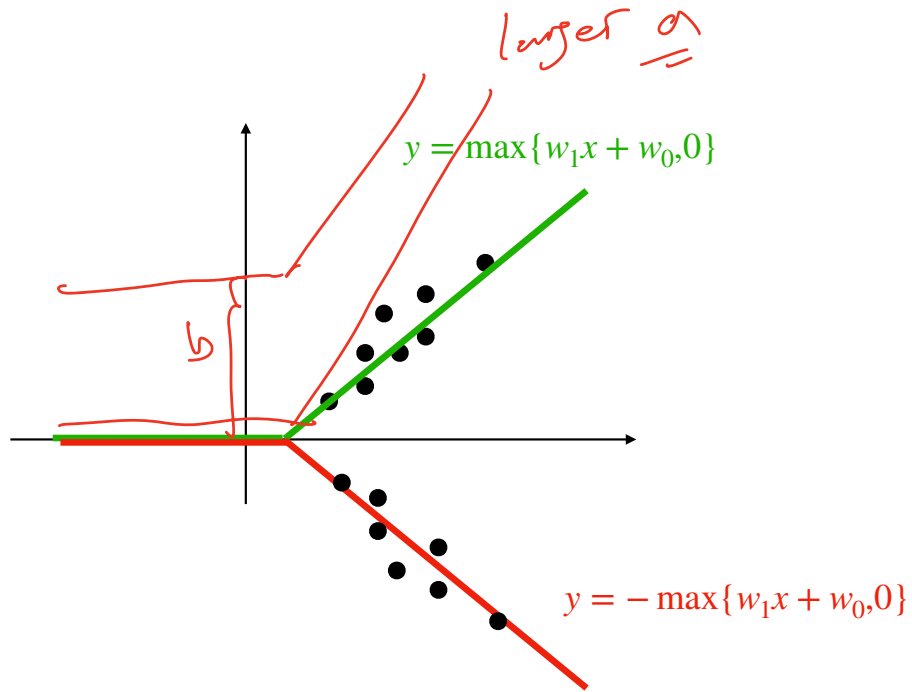


# A single neuron network

$$y = \max\{-w_1x - w_0, 0\}$$



# A single neuron network



$$y = a \max\{w_1x + w_0, 0\} + b$$

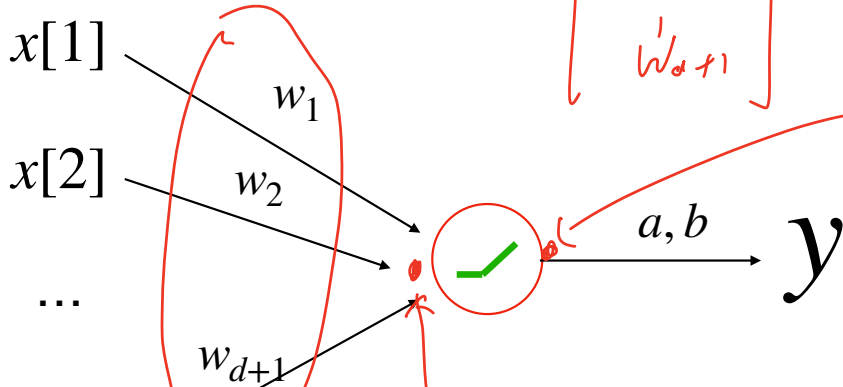
four parameters

$a, w_1, w_0, b$

# A single neuron network

$w_i \in \mathbb{R}$

$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_{d+1} \end{bmatrix}$$



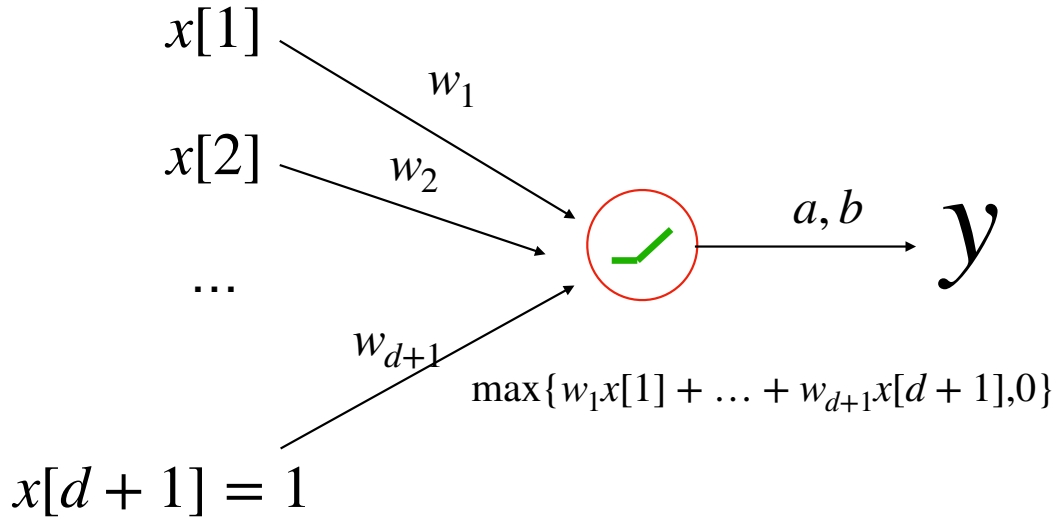
$$\max\{w^T x, 0\}$$

$$y = a \max\{w^T x, 0\} + b$$

$x[d+1] = 1$

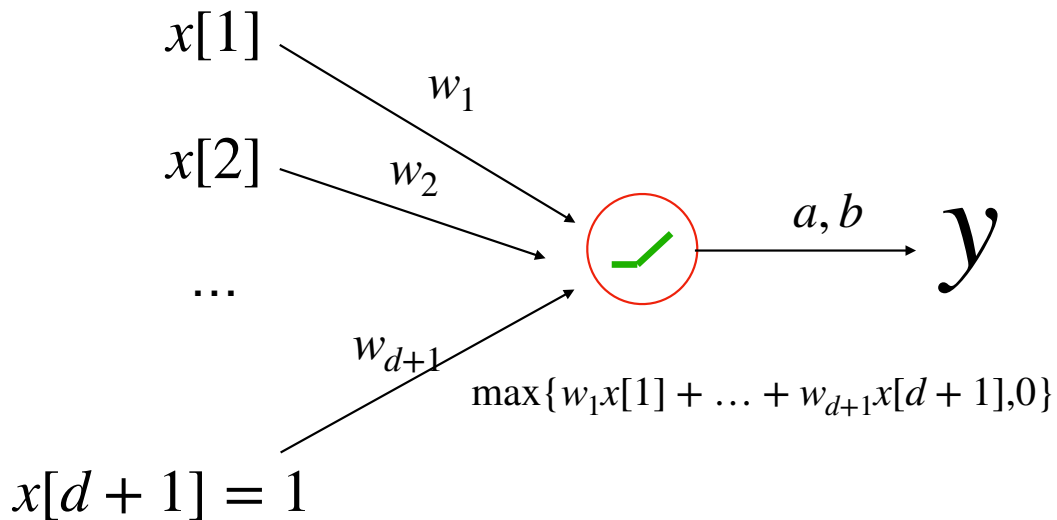
$$\sum_{i=1}^{d+1} w_i x[i] = w^T x$$

# A single neuron network

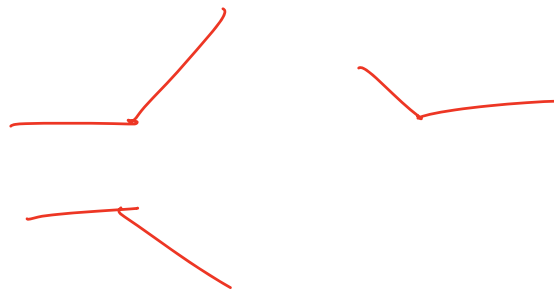


# A single neuron network

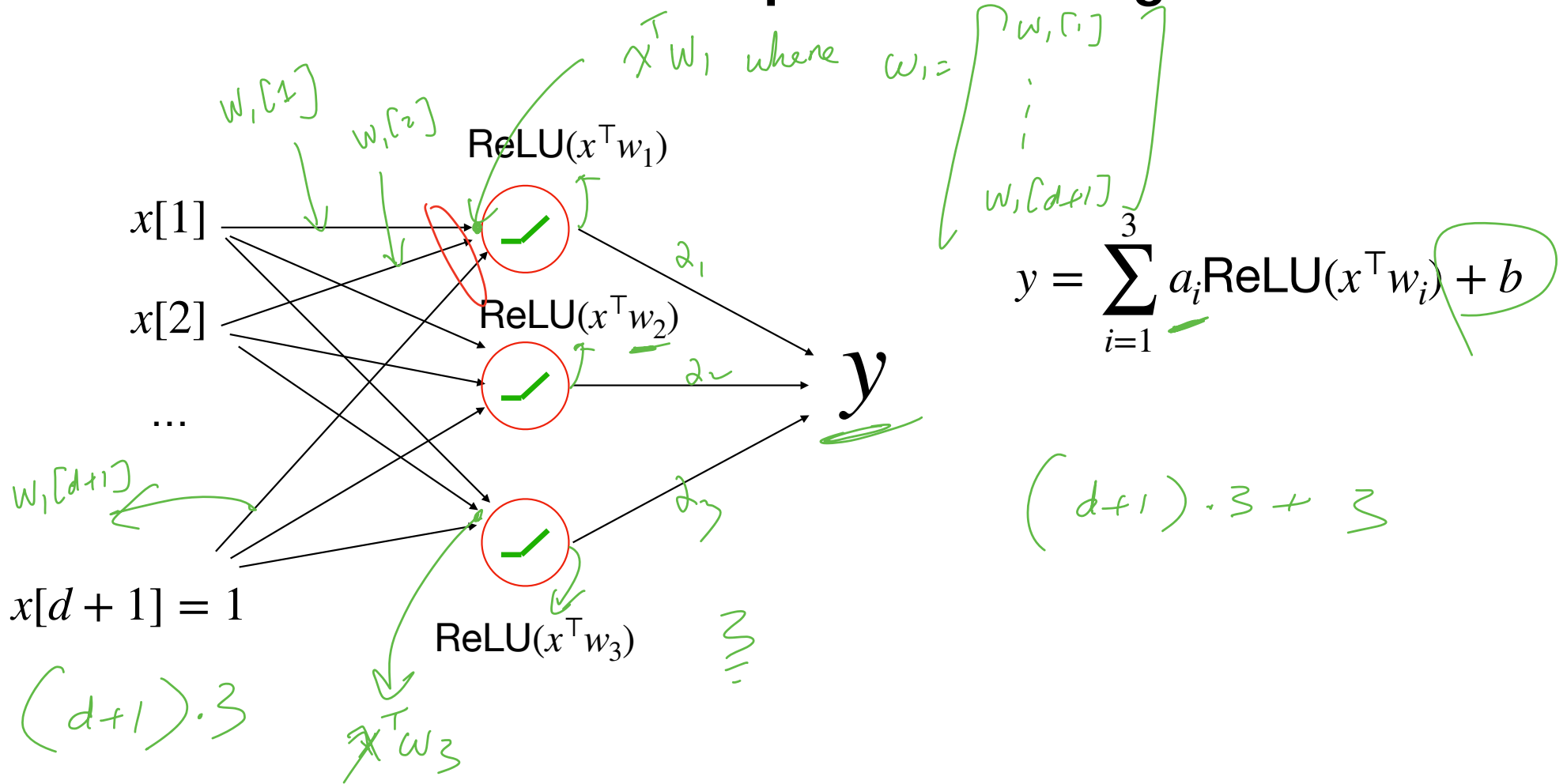
$$\max\{x, 0\} \Leftarrow \underline{\text{ReLU}(x)}$$



$$y = a\text{ReLU}(w^T x) + b$$

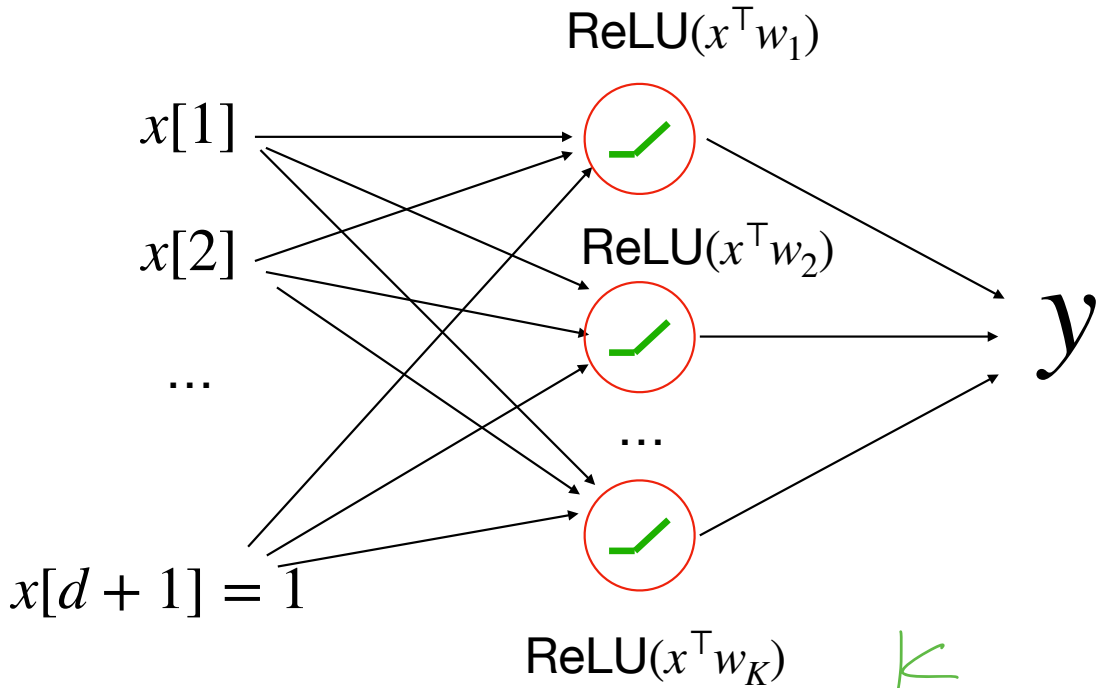


# Let us stack multiple neurons together



# Let us stack multiple neurons together

$K$  # of neurons

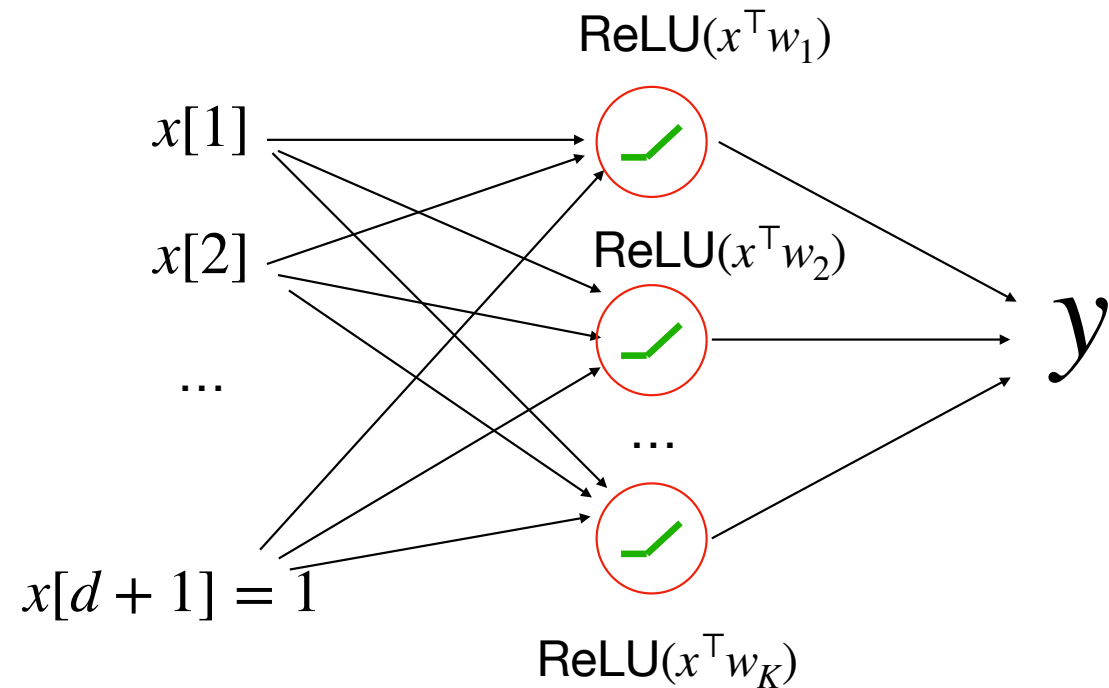


$(d+1) \cdot K$

$\text{ReLU}(x^T w_K)$   $K$



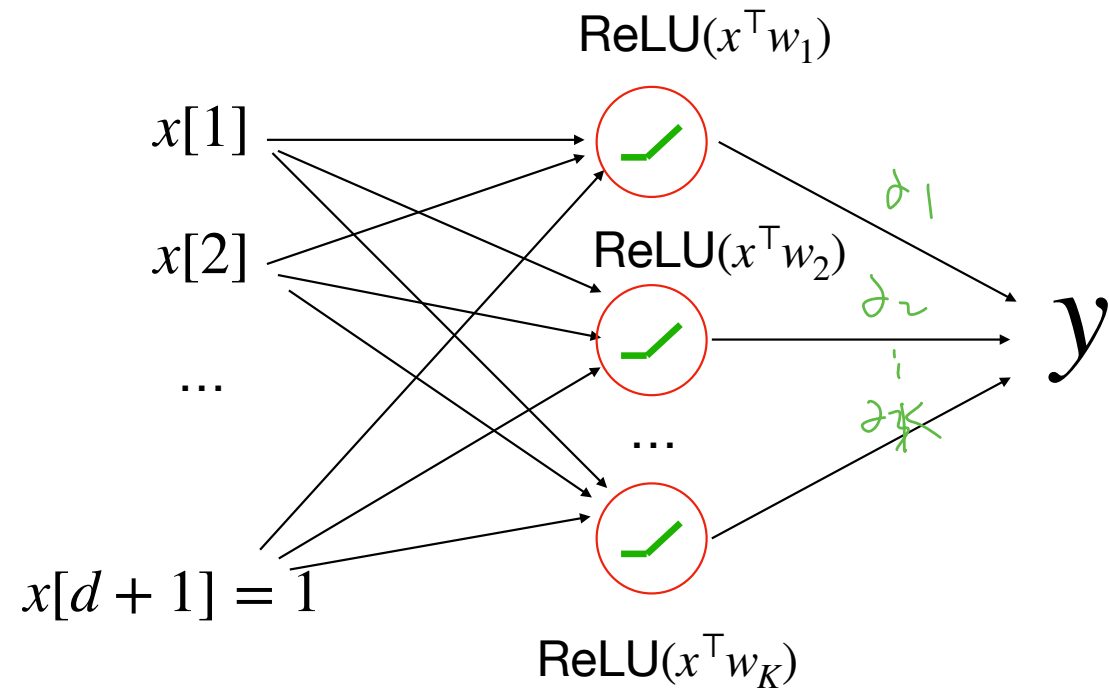
# Let us stack multiple neurons together



Vectorized form:

Define  $W = \begin{bmatrix} (w_1)^\top \\ \dots \\ (w_K)^\top \end{bmatrix} \in \mathbb{R}^{K \times (d+1)}$

# Let us stack multiple neurons together



Vectorized form:

$$\text{Define } W = \begin{bmatrix} (w_1)^\top \\ \dots \\ (w_K)^\top \end{bmatrix} \in \mathbb{R}^{K \times (d+1)}$$

$$\alpha = [a_1, \dots, a_K]^\top$$

# Let us stack multiple neurons together

$$\text{ReLU}(z)$$

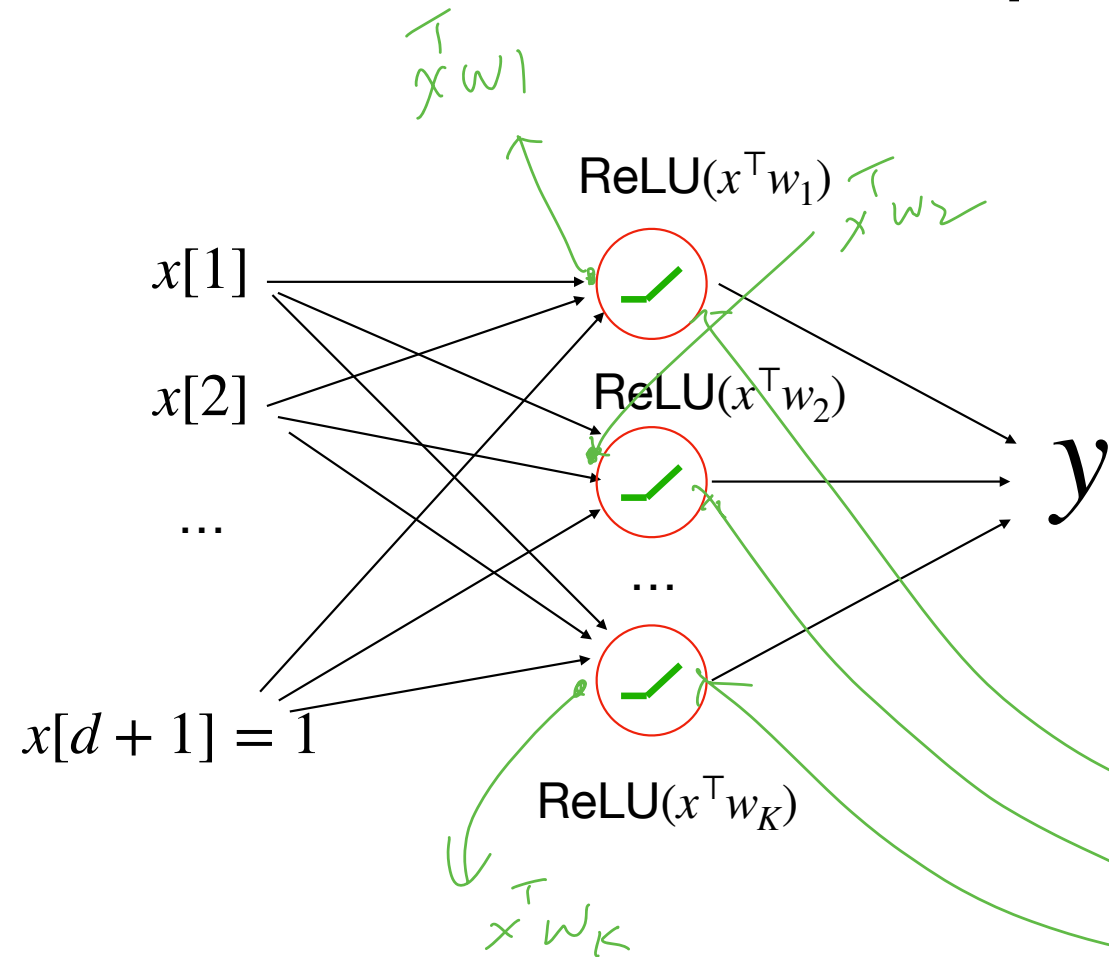
$$= \begin{bmatrix} \text{ReLU}(z_1) \\ \vdots \\ \text{ReLU}(z_d) \end{bmatrix}$$

Vectorized form:

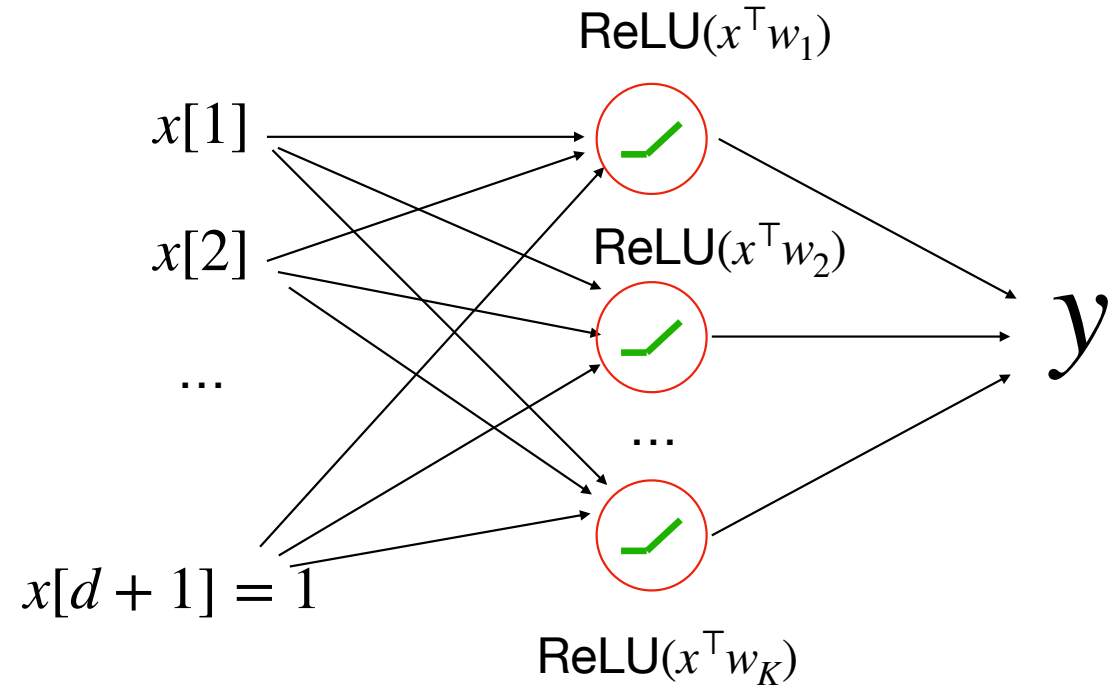
$$\text{Define } W = \begin{bmatrix} (w_1)^T \\ \dots \\ (w_K)^T \end{bmatrix} \in \mathbb{R}^{K \times (d+1)}$$

$$\alpha = [a_1, \dots, a_K]^T$$

$$y = \alpha^T (\text{ReLU}(Wx)) + b$$



# Let us stack multiple neurons together



Vectorized form:

Define  $W = \begin{bmatrix} (w_1)^\top \\ \dots \\ (w_K)^\top \end{bmatrix} \in \mathbb{R}^{K \times (d+1)}$

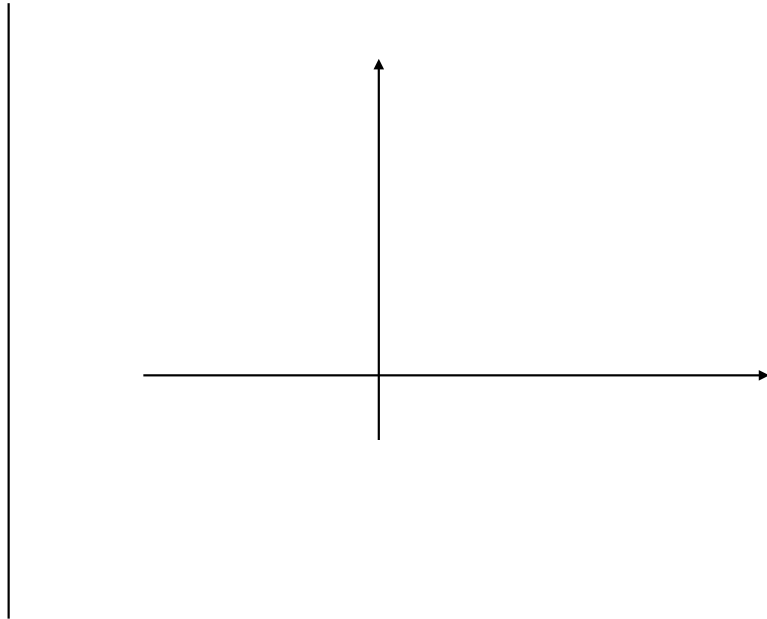
$$\alpha = [a_1, \dots, a_K]^\top$$

$$y = \alpha^\top (\text{ReLU}(Wx)) + b$$

Learable feature  $\phi(x)$

# What does a neural network approximate

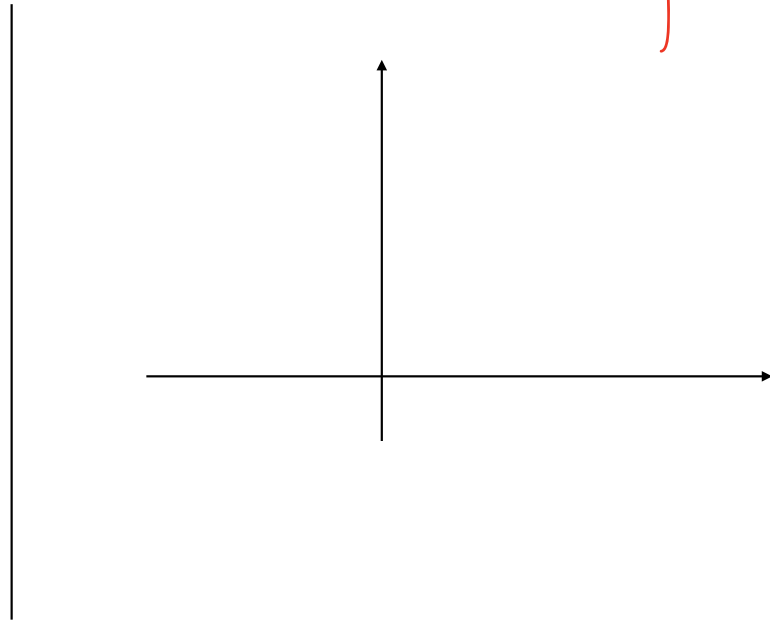
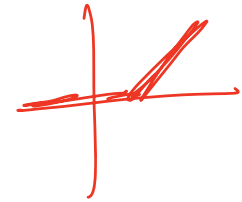
$$y = \alpha^T (\text{ReLU}(Wx)) + b$$



# What does a neural network approximate

$$y = \alpha^T (\text{ReLU}(Wx)) + b = \sum_{i=1}^K \alpha_i \text{ReLU}(w_i^T x) + b$$

It's a pieces wise linear functions

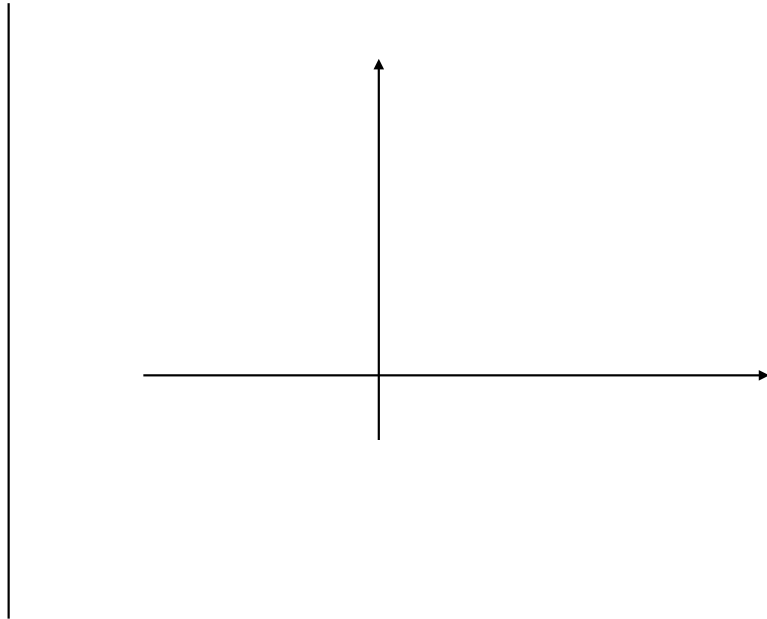


# What does a neural network approximate

$$y = \alpha^T (\text{ReLU}(Wx)) + b$$

It's a pieces wise linear functions

Consider  $d = 1$  case (and assume  $b = 0$ ):



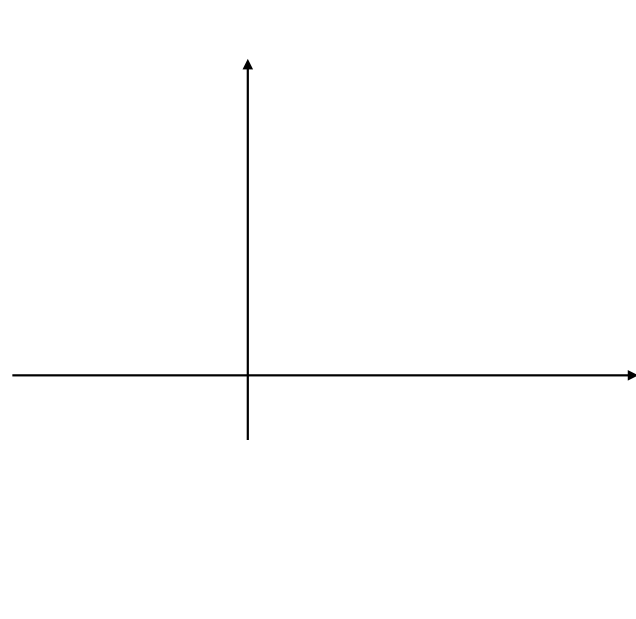
# What does a neural network approximate

$$y = \alpha^\top (\text{ReLU}(Wx)) + b$$

It's a pieces wise linear functions

Consider  $d = 1$  case (and assume  $b = 0$ ):

$$K = 1 : y = a_1 \max\{w_1 x + c_1, 0\}$$





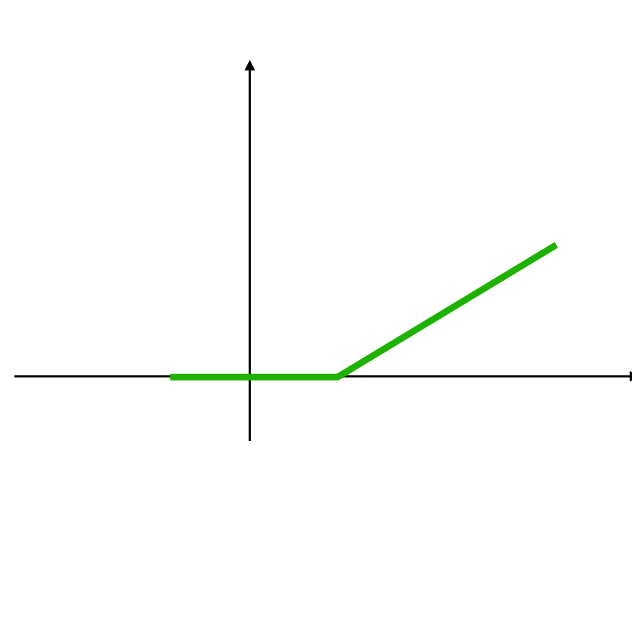
# What does a neural network approximate

$$y = \alpha^\top (\text{ReLU}(Wx)) + b$$

It's a pieces wise linear functions

Consider  $d = 1$  case (and assume  $b = 0$ ):

$$K = 1 : y = a_1 \max\{w_1 x + c_1, 0\}$$



# What does a neural network approximate

$$y = \alpha^T (\text{ReLU}(Wx)) + b$$

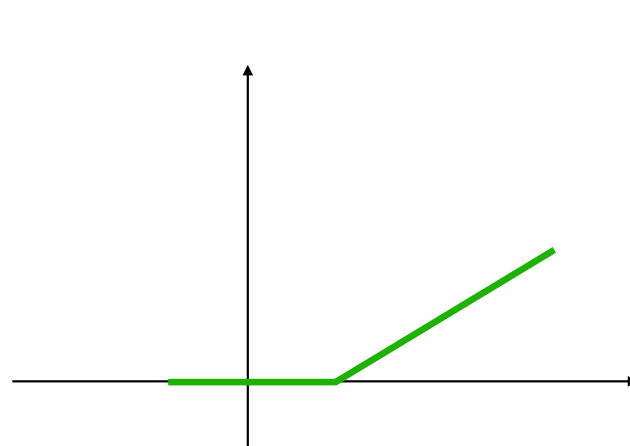
It's a pieces wise linear functions

Consider  $d = 1$  case (and assume  $b = 0$ ):

$$K = 1 : y = a_1 \max\{w_1x + c_1, 0\}$$

$$K = 2 : y = a_1 \max\{w_1x + c_1, 0\}$$

$$+ \underline{a_2 \max\{w_2x + c_2, 0\}}$$



# What does a neural network approximate

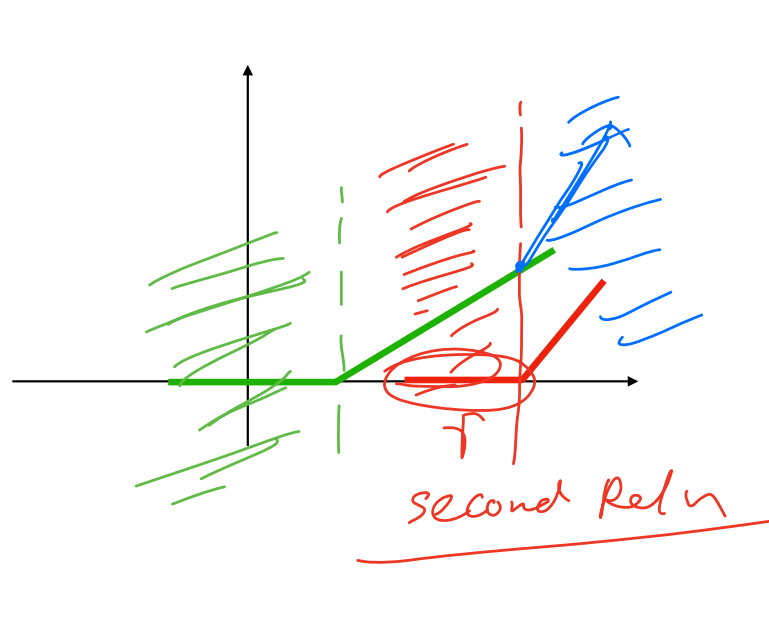
$$y = \alpha^T (\text{ReLU}(Wx)) + b$$

It's a pieces wise linear functions

Consider  $d = 1$  case (and assume  $b = 0$ ):

$$K = 1 : y = a_1 \max\{w_1x + c_1, 0\}$$

$$K = 2 : y = a_1 \max\{w_1x + c_1, 0\} \\ + a_2 \max\{w_2x + c_2, 0\}$$



# What does a neural network approximate

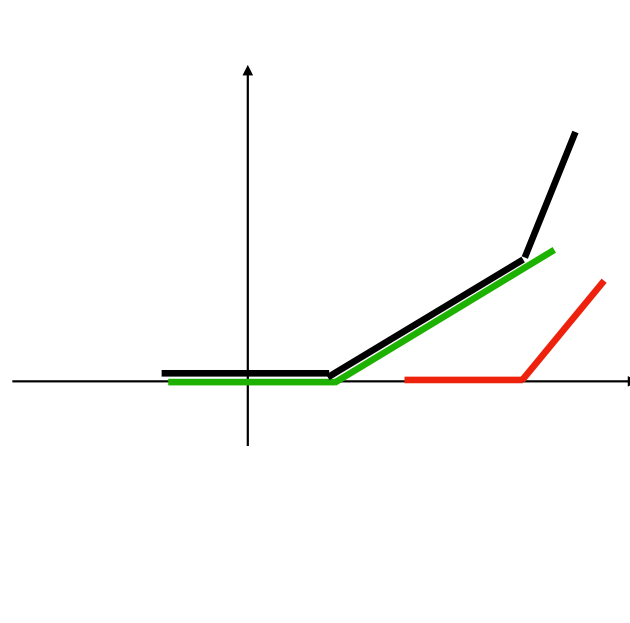
$$y = \alpha^T (\text{ReLU}(Wx)) + b$$

It's a pieces wise linear functions

Consider  $d = 1$  case (and assume  $b = 0$ ):

$$K = 1 : y = a_1 \max\{w_1x + c_1, 0\}$$

$$K = 2 : y = a_1 \max\{w_1x + c_1, 0\} \\ + a_2 \max\{w_2x + c_2, 0\}$$



# What does a neural network approximate

$$y = \alpha^T (\text{ReLU}(Wx)) + b$$

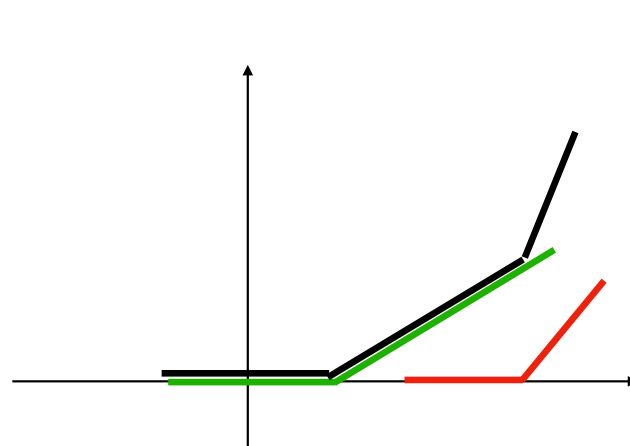
It's a pieces wise linear functions

Consider  $d = 1$  case (and assume  $b = 0$ ):

$$K = 1 : y = a_1 \max\{w_1x + c_1, 0\}$$

$$K = 2 : y = a_1 \max\{w_1x + c_1, 0\} \\ + a_2 \max\{w_2x + c_2, 0\}$$

$$K = 3 : y = a_1 \max\{w_1x + c_1, 0\} \\ + a_2 \max\{w_2x + c_2, 0\} \\ + a_3 \max\{w_3x + c_3, 0\}$$



# What does a neural network approximate

$$y = \alpha^T (\text{ReLU}(Wx)) + b$$

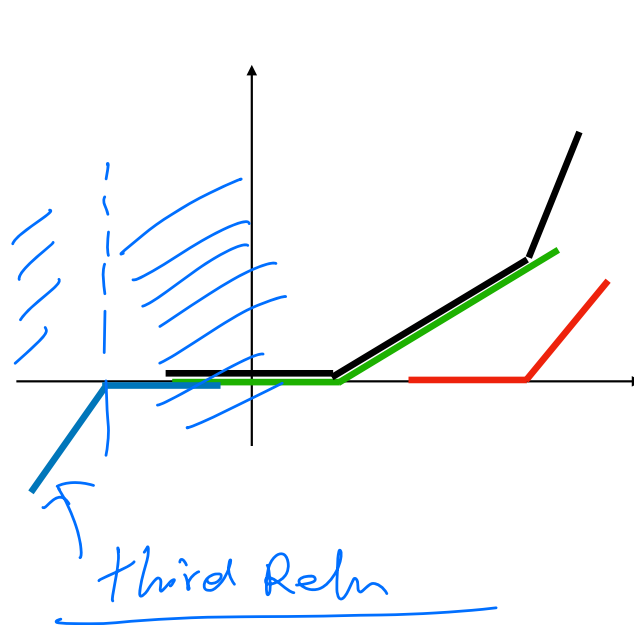
It's a pieces wise linear functions

Consider  $d = 1$  case (and assume  $b = 0$ ):

$$K = 1 : y = a_1 \max\{w_1x + c_1, 0\}$$

$$K = 2 : y = a_1 \max\{w_1x + c_1, 0\} \\ + a_2 \max\{w_2x + c_2, 0\}$$

$$K = 3 : y = a_1 \max\{w_1x + c_1, 0\} \\ + a_2 \max\{w_2x + c_2, 0\} \\ + a_3 \max\{w_3x + c_3, 0\}$$



# What does a neural network approximate

$$y = \alpha^\top (\text{ReLU}(Wx)) + b$$

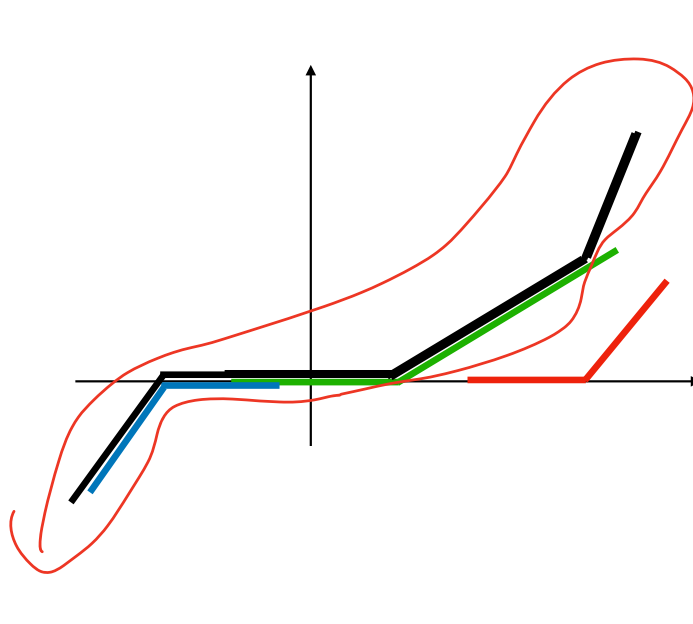
It's a pieces wise linear functions

Consider  $d = 1$  case (and assume  $b = 0$ ):

$$K = 1 : y = a_1 \max\{w_1x + c_1, 0\}$$

$$K = 2 : y = a_1 \max\{w_1x + c_1, 0\} \\ + a_2 \max\{w_2x + c_2, 0\}$$

$$K = 3 : y = a_1 \max\{w_1x + c_1, 0\} \\ + a_2 \max\{w_2x + c_2, 0\} \\ + a_3 \max\{w_3x + c_3, 0\}$$



# What does a neural network approximate

$$y = \alpha^T (\text{ReLU}(Wx)) + b$$

Claim: a wide enough one layer NN can approximate any smooth functions

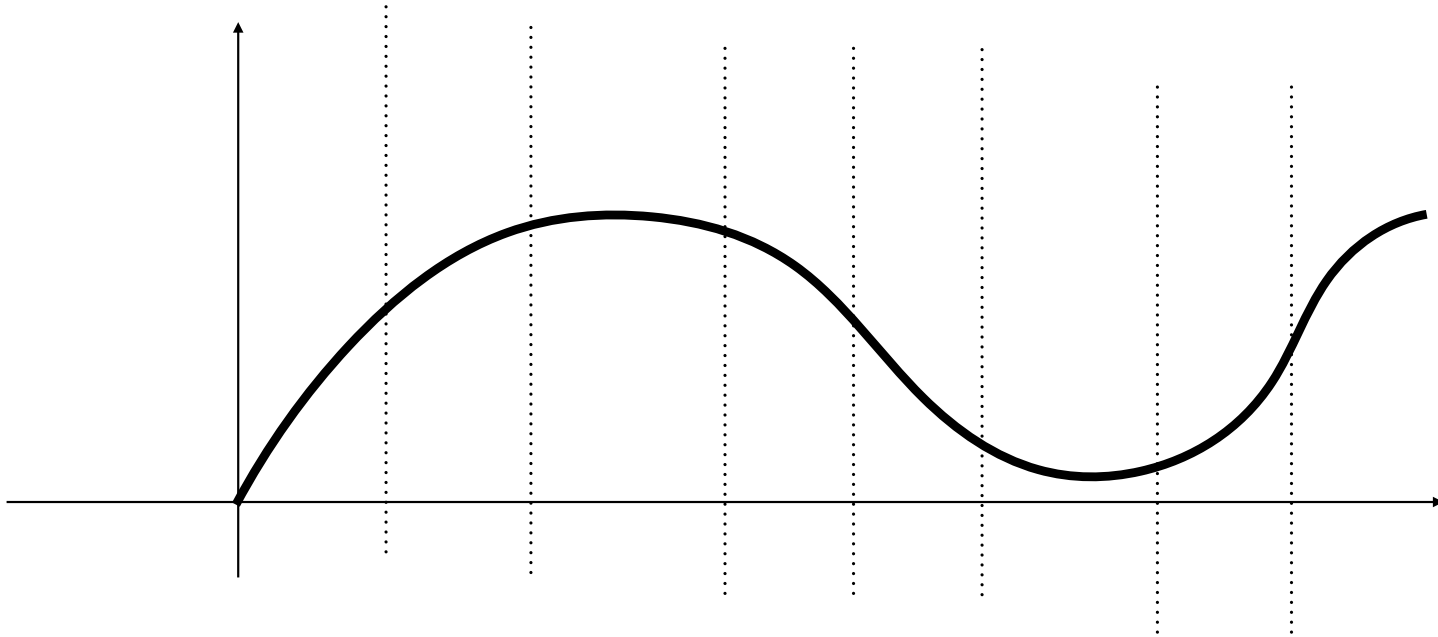




# What does a neural network approximate

$$y = \alpha^T (\text{ReLU}(Wx)) + b$$

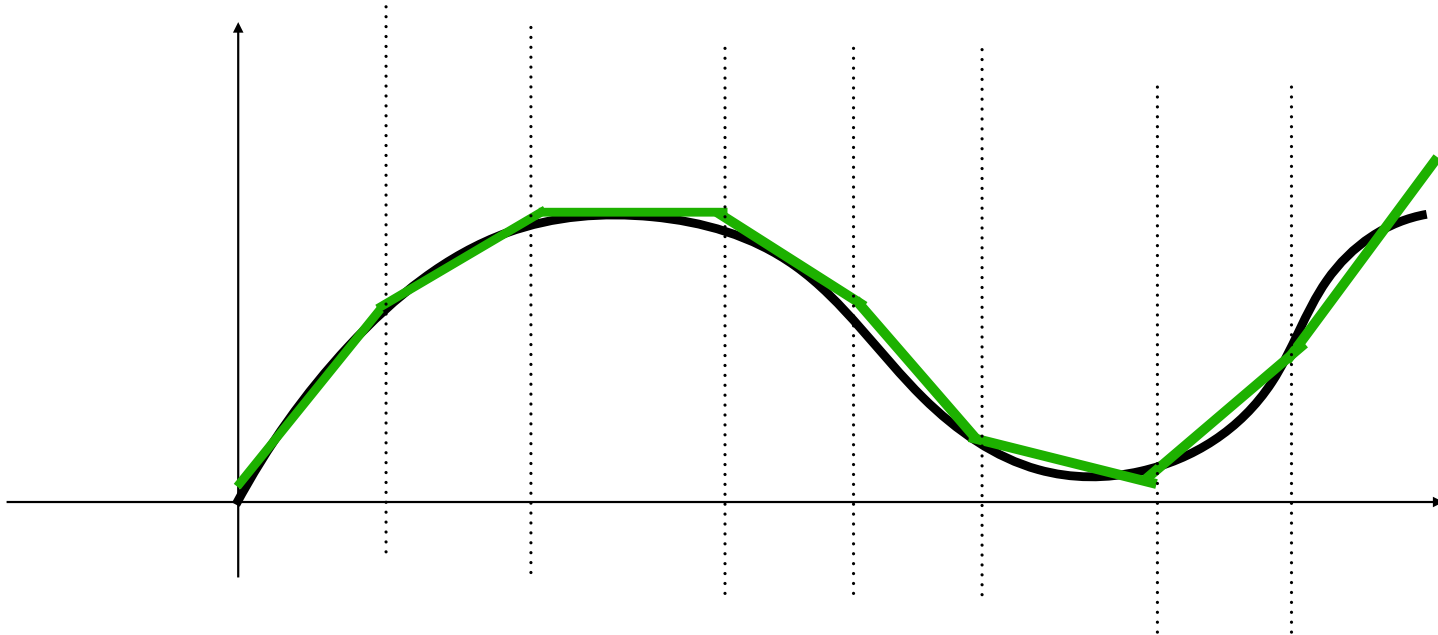
Claim: a wide enough one layer NN can approximate any smooth functions



# What does a neural network approximate

$$y = \alpha^T (\text{ReLU}(Wx)) + b$$

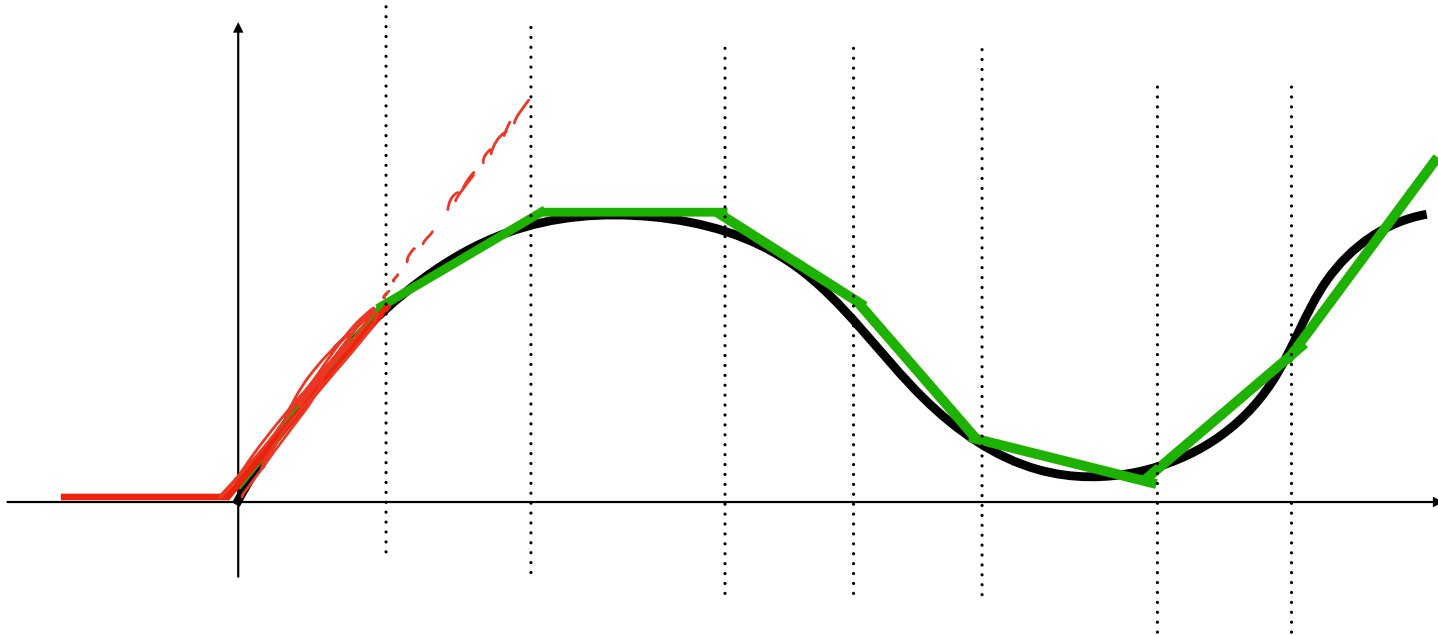
Claim: a wide enough one layer NN can approximate any smooth functions



# What does a neural network approximate

$$y = \alpha^T (\text{ReLU}(Wx)) + b$$

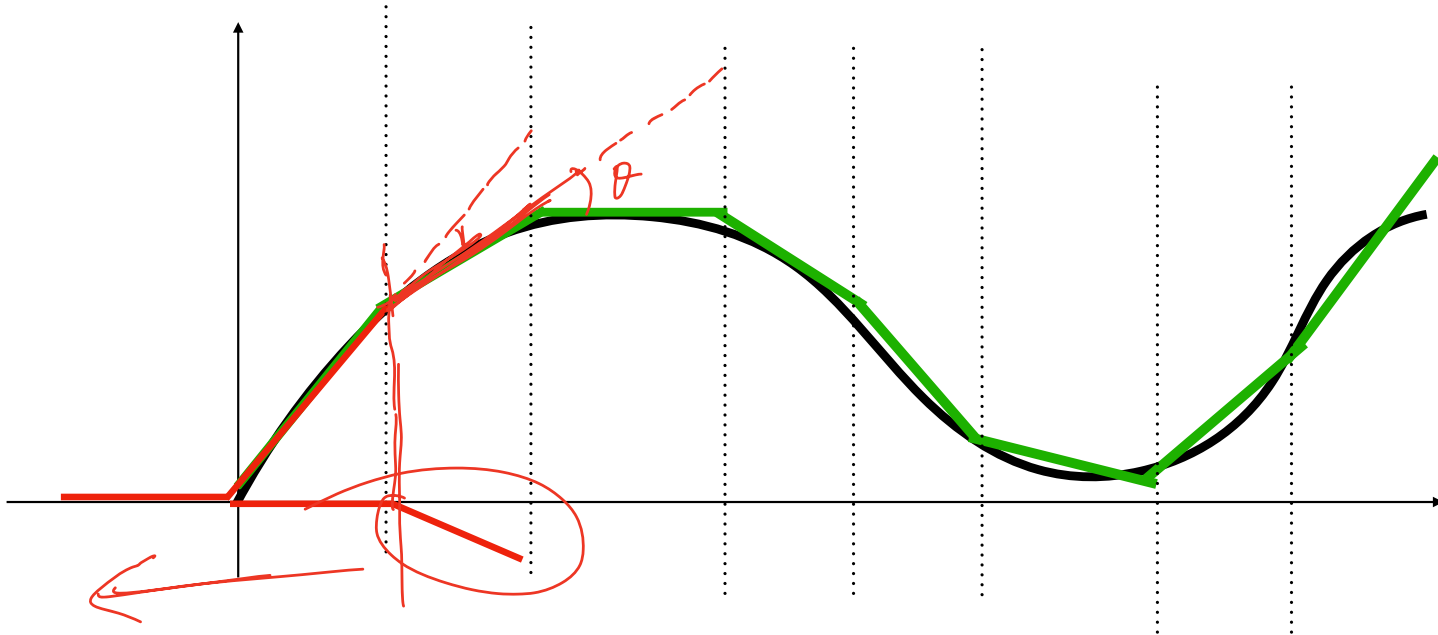
Claim: a wide enough one layer NN can approximate any smooth functions



# What does a neural network approximate

$$y = \alpha^T (\text{ReLU}(Wx)) + b$$

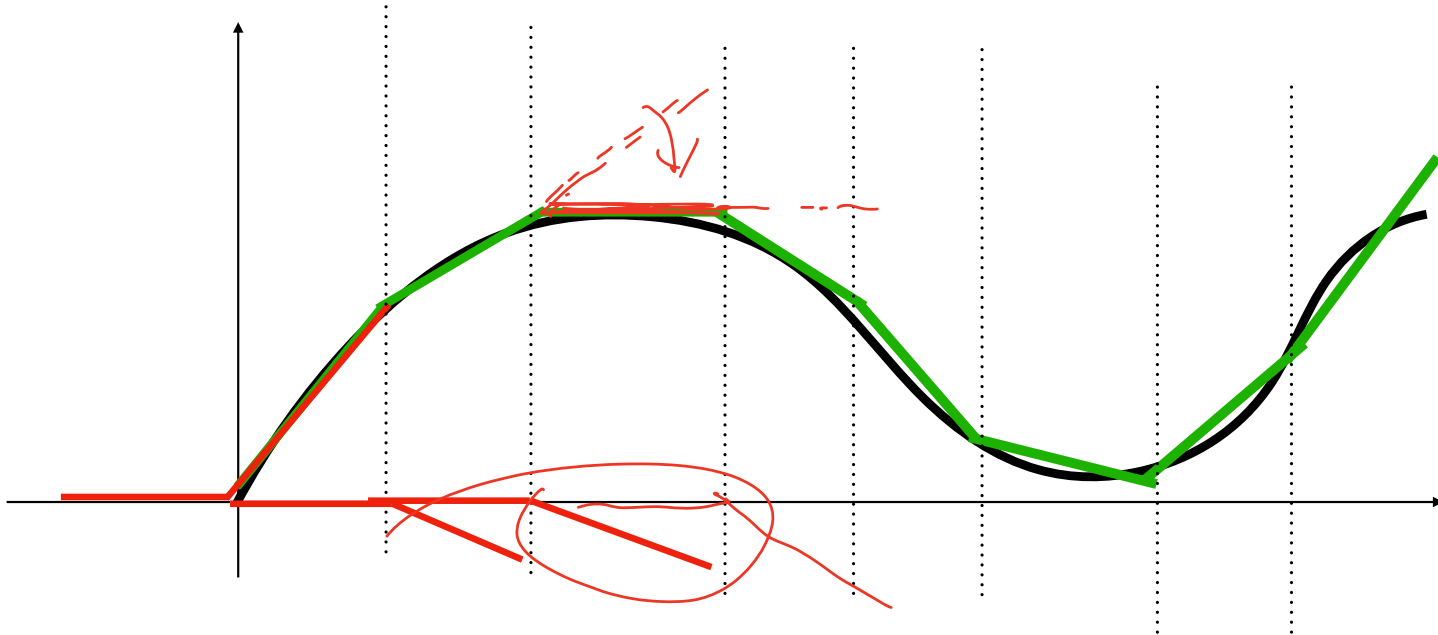
Claim: a wide enough one layer NN can approximate any smooth functions



# What does a neural network approximate

$$y = \alpha^T (\text{ReLU}(Wx)) + b$$

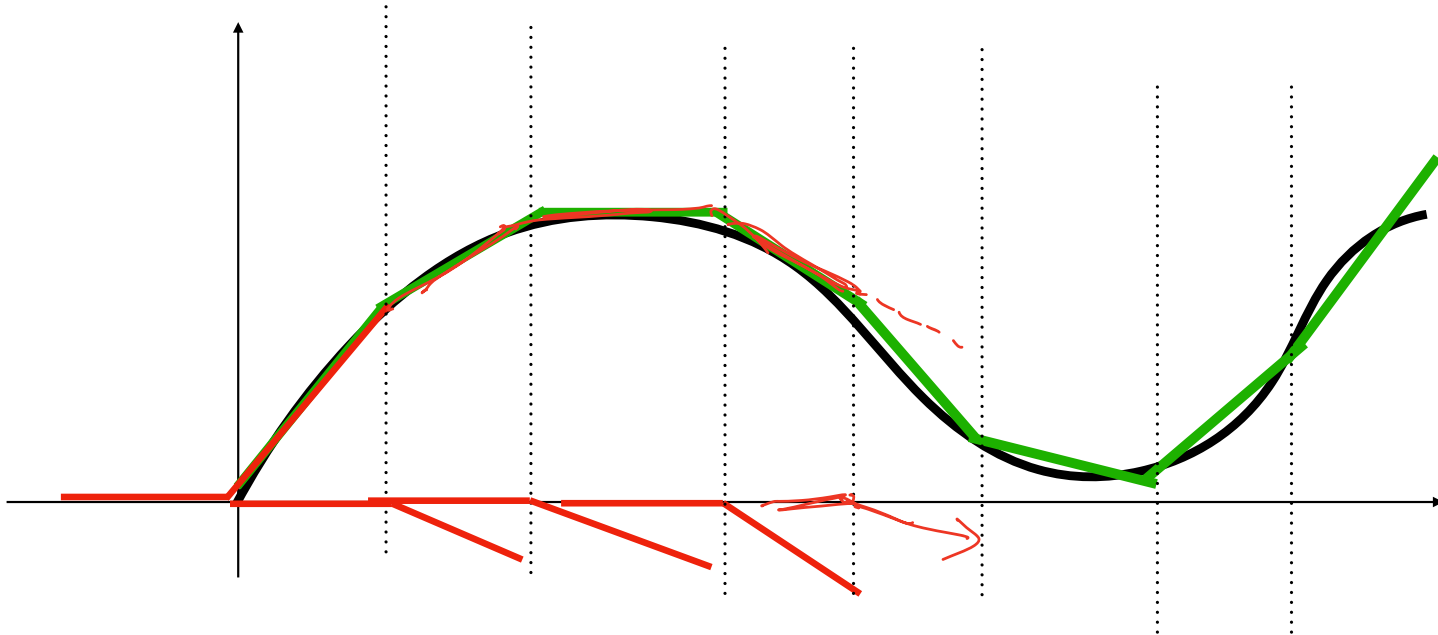
Claim: a wide enough one layer NN can approximate any smooth functions



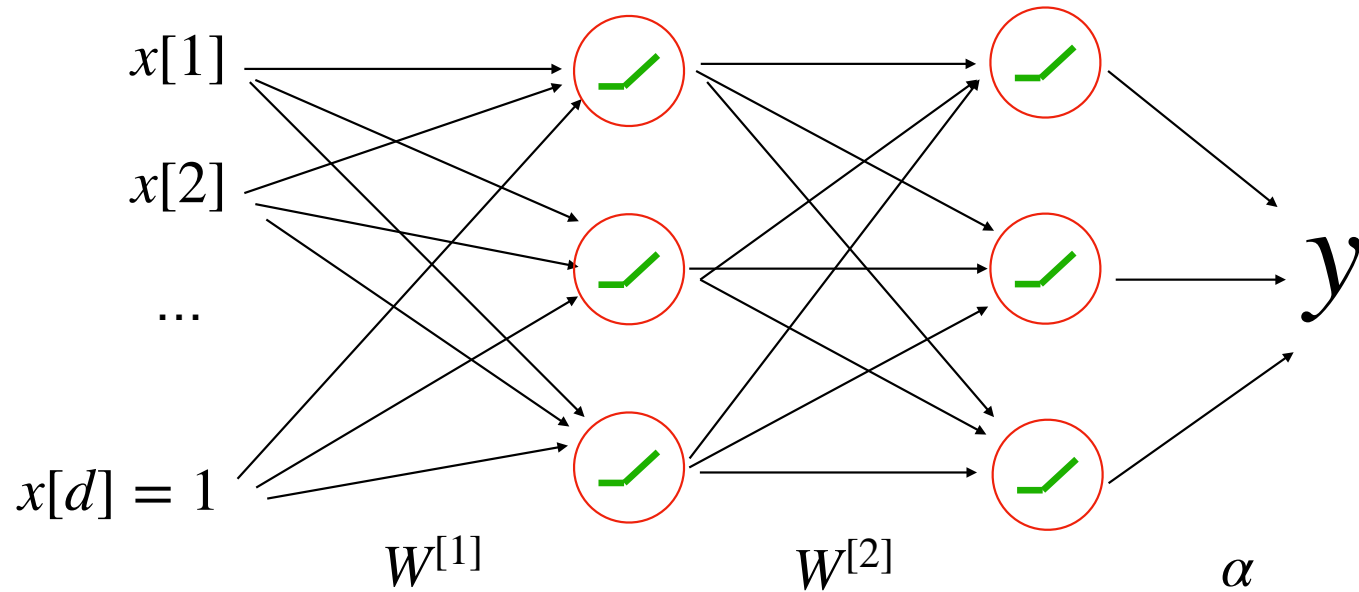
# What does a neural network approximate

$$y = \alpha^T (\text{ReLU}(Wx)) + b$$

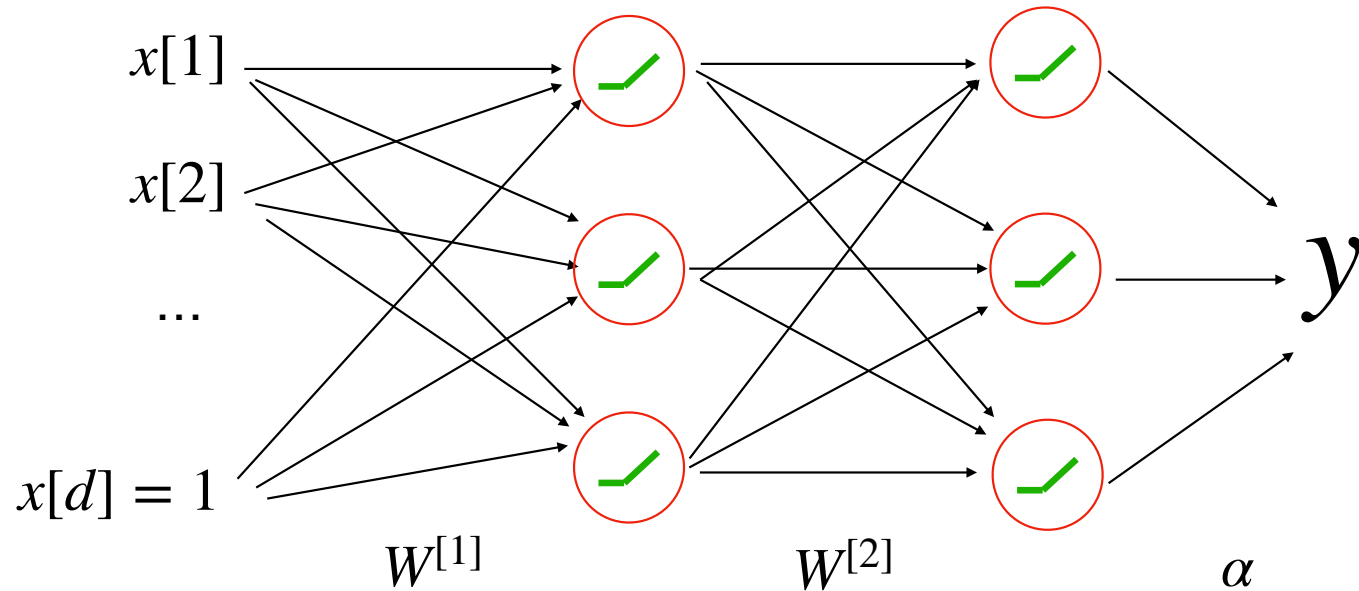
Claim: a wide enough one layer NN can approximate any smooth functions



# A multi-layer fully connected neural network



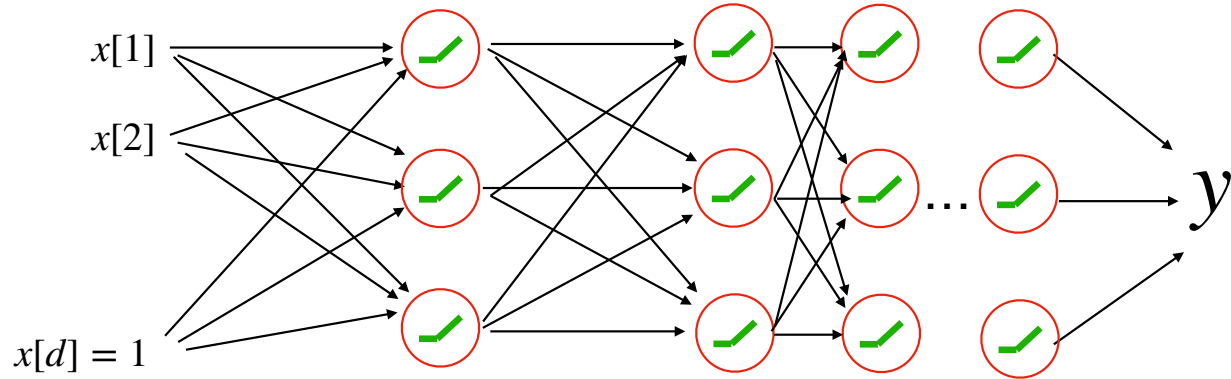
# A multi-layer fully connected neural network



$$y = \alpha^T \text{ReLU} \left( W^{[2]} \text{ReLU} \left( W^{[1]} x \right) \right) + b$$

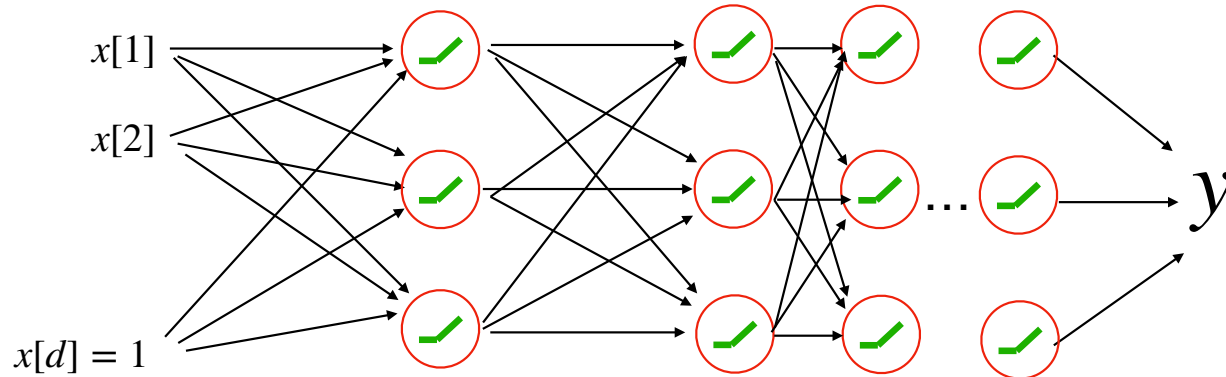


# A multi-layer fully connected neural network



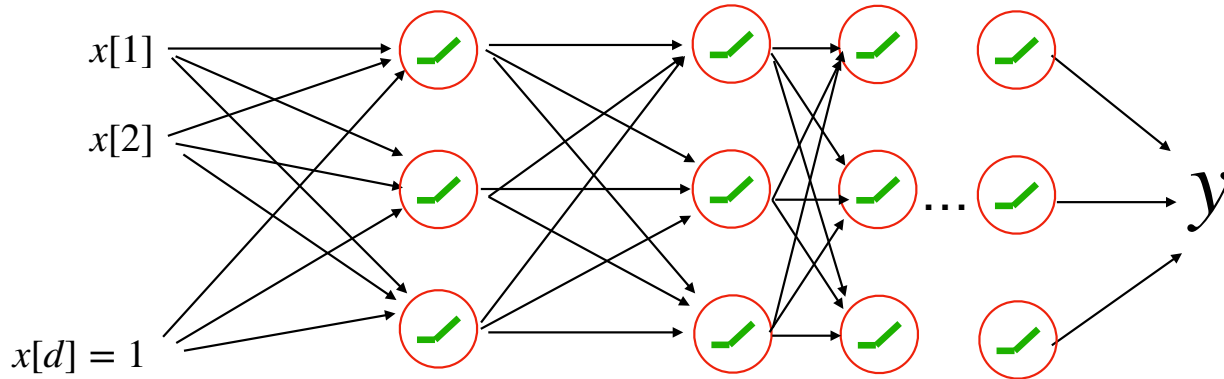
# A multi-layer fully connected neural network

Define it by a forward pass:



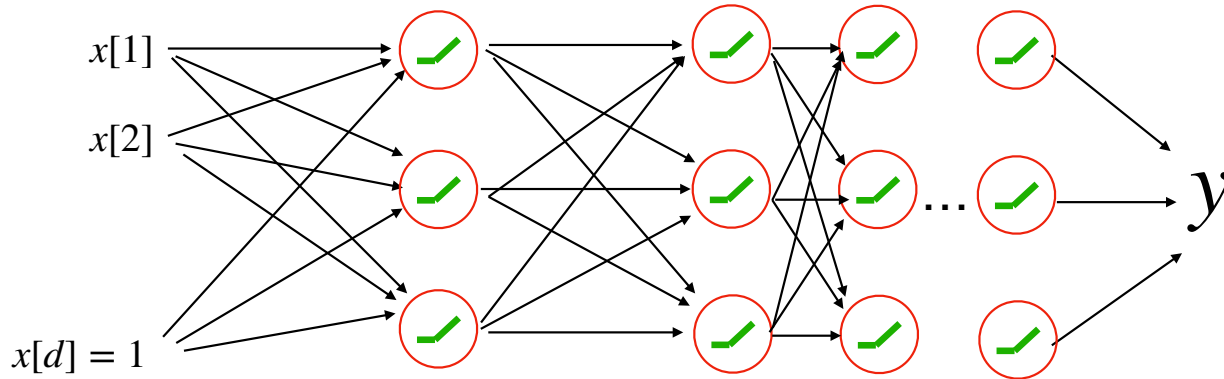
# A multi-layer fully connected neural network

Define it by a forward pass:



$$z^{[1]} = x$$

# A multi-layer fully connected neural network

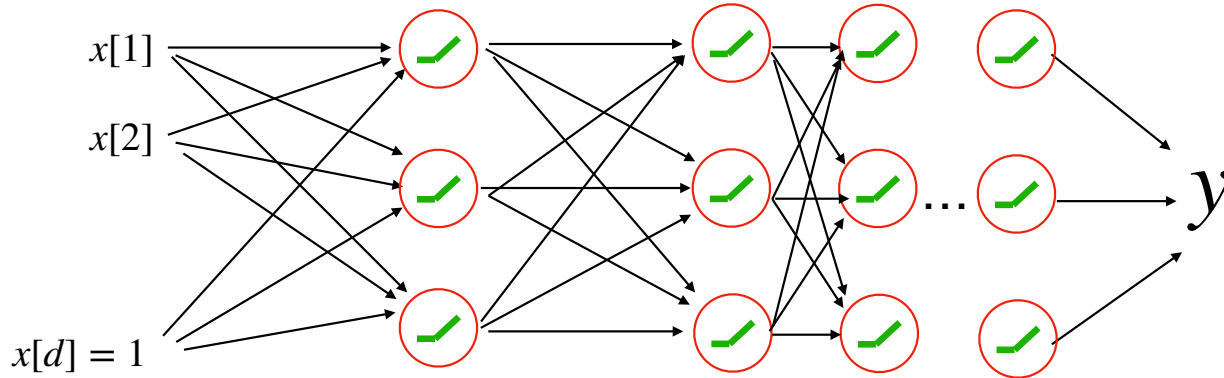


Define it by a forward pass:

$$z^{[1]} = x$$

For  $t = 1$  to  $T-1$ :

# A multi-layer fully connected neural network



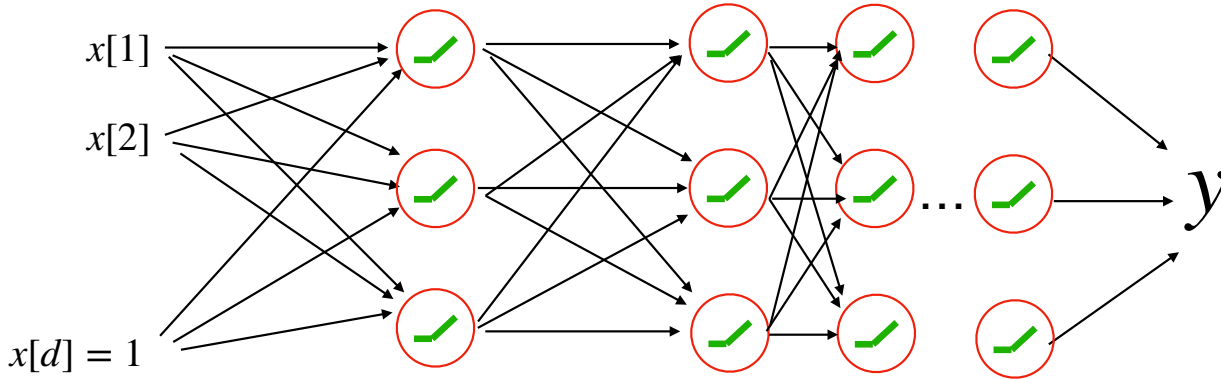
Define it by a forward pass:

$$z^{[1]} = x$$

For  $t = 1$  to  $T-1$ :

$$z^{[t+1]} = \text{ReLU} (W^{[t]}z^t)$$

# A multi-layer fully connected neural network



Define it by a forward pass:

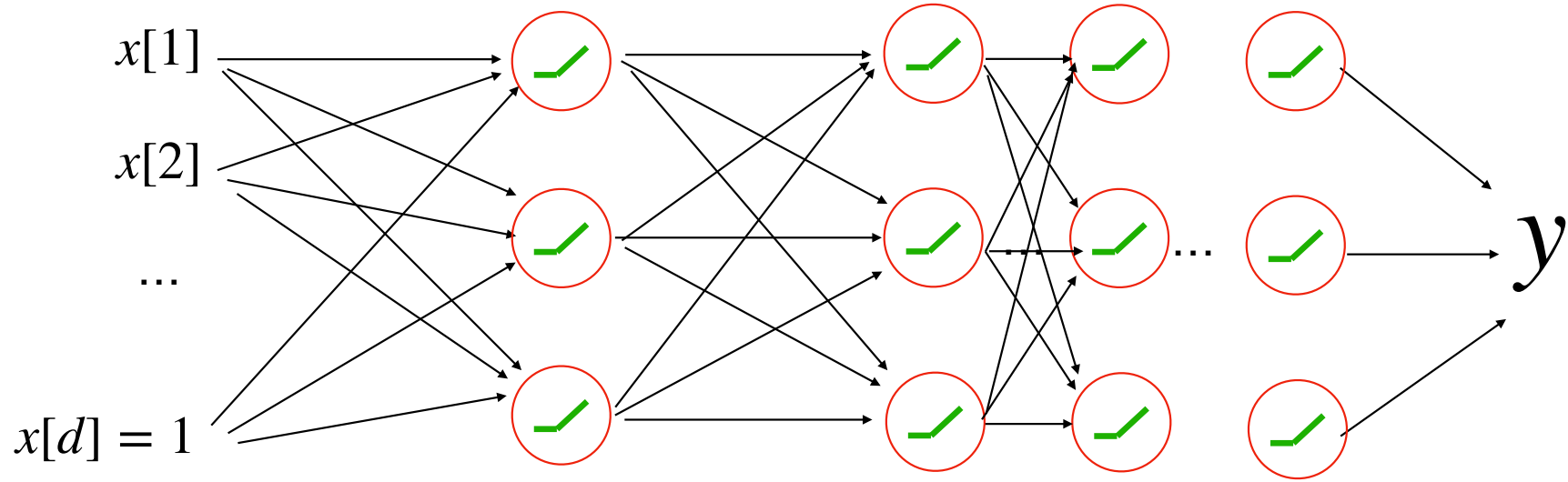
$$z^{[1]} = x$$

For  $t = 1$  to  $T-1$ :

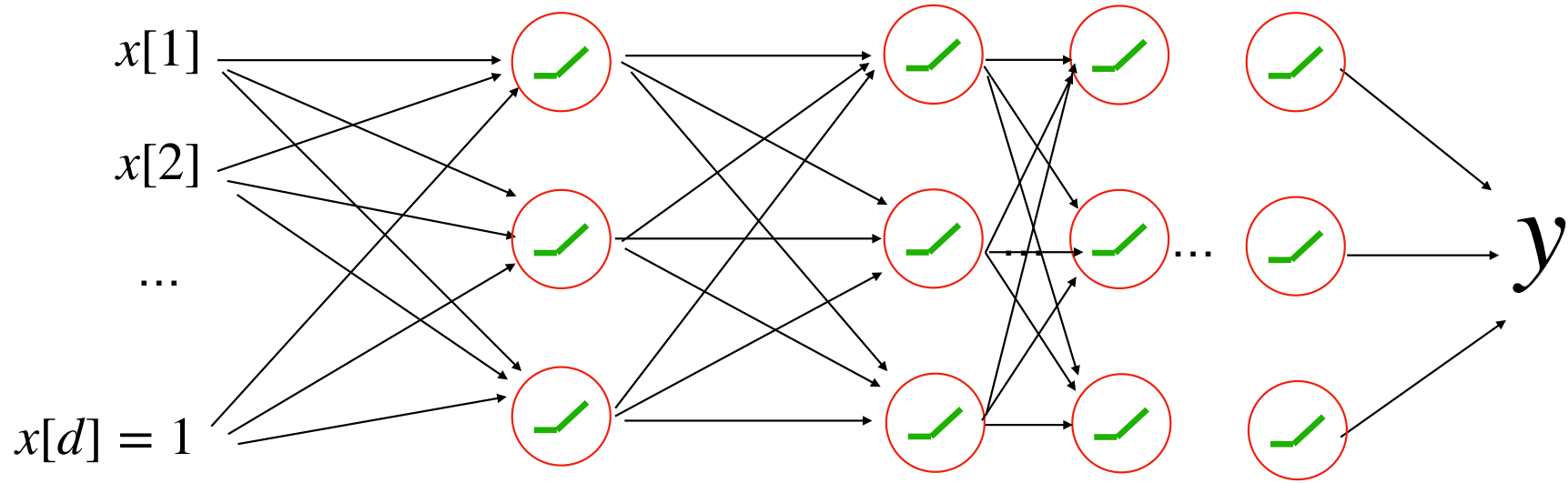
$$z^{[t+1]} = \text{ReLU} (W^{[t]} z^t)$$

$$y = \alpha^\top z^{[T]} + b$$

# The benefits of going deep



# The benefits of going deep



Allows us to represent complicated functions without making NN too wide



# Summary for today

Neural network is universal function approximation

Next lecture: backpropagation for computing gradients efficiently