

## Class Plan

Reminders

Recap

Neural Nets

Architectures

Gradient Descent

Back propagation

Grading 100 points

93 - A 86 - B+

90 - A- 83 - B

## Google Form

- go/cs101/
- due midnight every class day
- worth 1 point per submission

## Homework

- 4 points (3 soln., 1 self grade)
- use LaTeX for theory
- python notebook for code

24 hours no questions asked  
total lateness

## Honor code

Attribute ideas and code

My expectation:

- do reading before class  
(max 15 min)
- do problem day it's  
assigned (mostly)
- complete google form

# Recap

3 step recipe

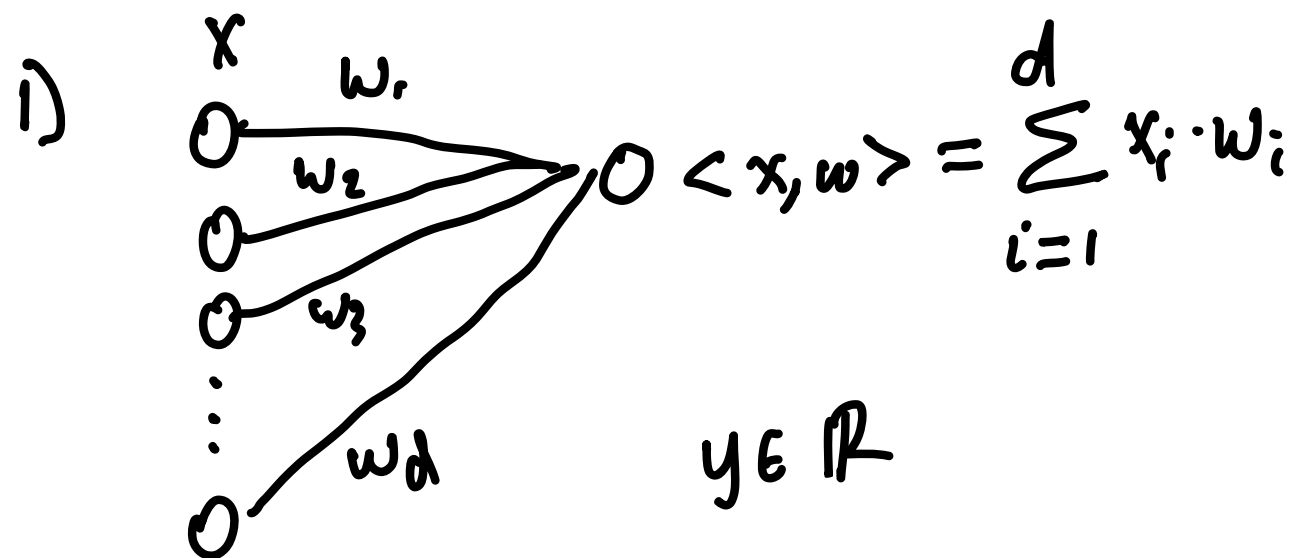
1) architecture / model

2) loss  $\frac{1}{n} \sum_{i=1}^n \ell(y_i, \hat{y}_i)$

3) optimizer

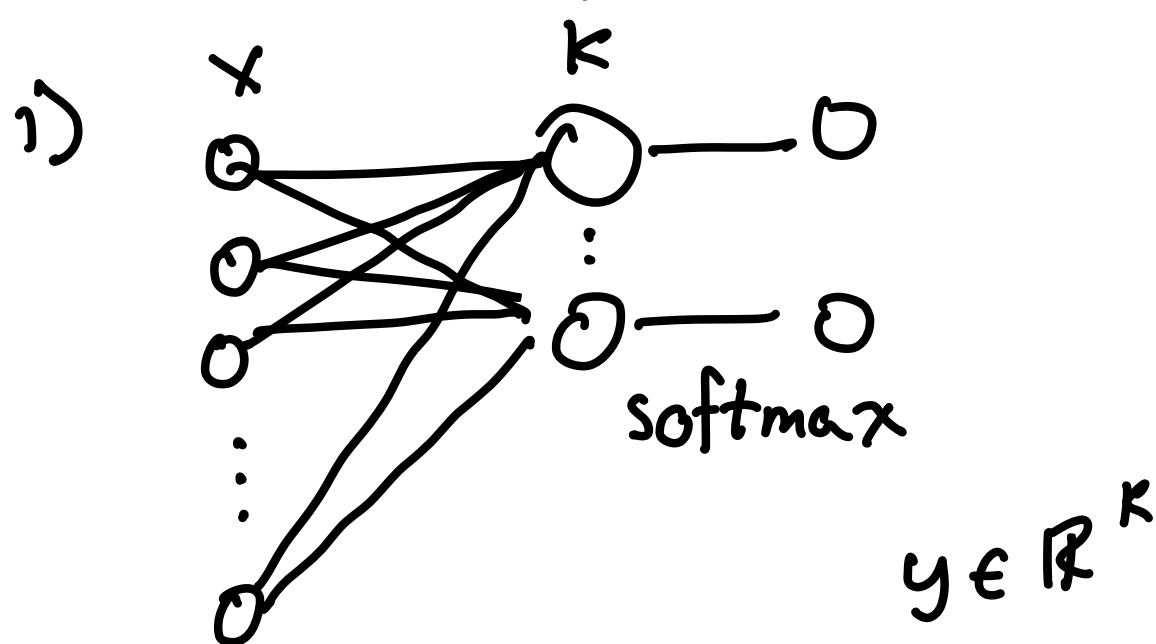
## Linear Regression

$$(x, y) \quad f_w(x) \approx y$$



$$2) \ell(y, \hat{y}) = \frac{1}{2} (y - \hat{y})^2$$

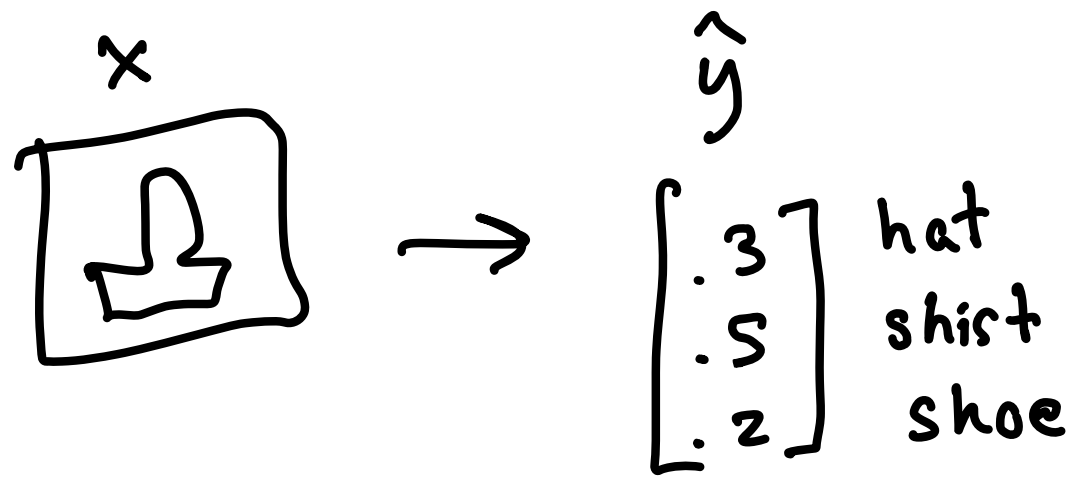
## Logistic Regression



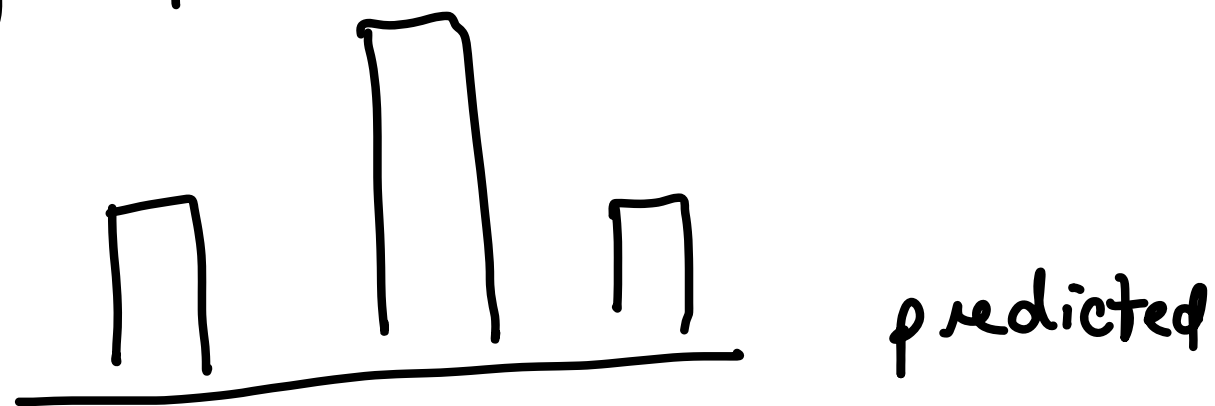
2)

$$\ell(y, \hat{y}) = \sum_{j=1}^k \mathbb{1}[y=j] \cdot -\log(\hat{y}_j)$$

$\uparrow$   
1 if  $y=j$   
0 else



$y = 1$

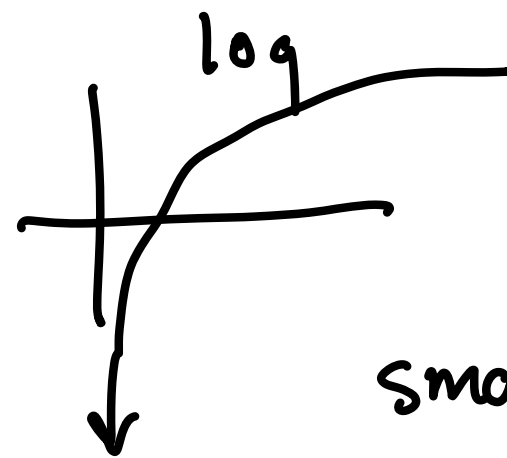


hat shirt shoe



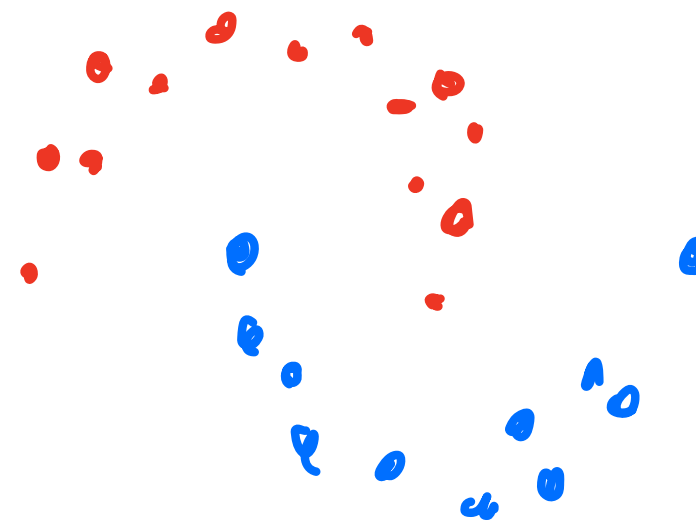
$$l(y, \hat{y}) = \sum_{j=1}^k \mathbb{1}[y=j] \cdot -\log(\hat{y}_j)$$

$\uparrow$   
 1 if  $y=j$   
 0 else



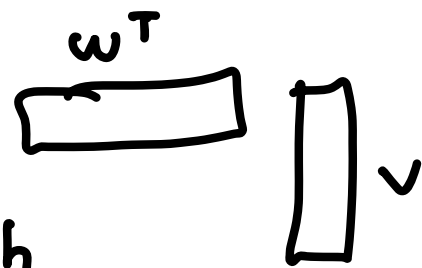
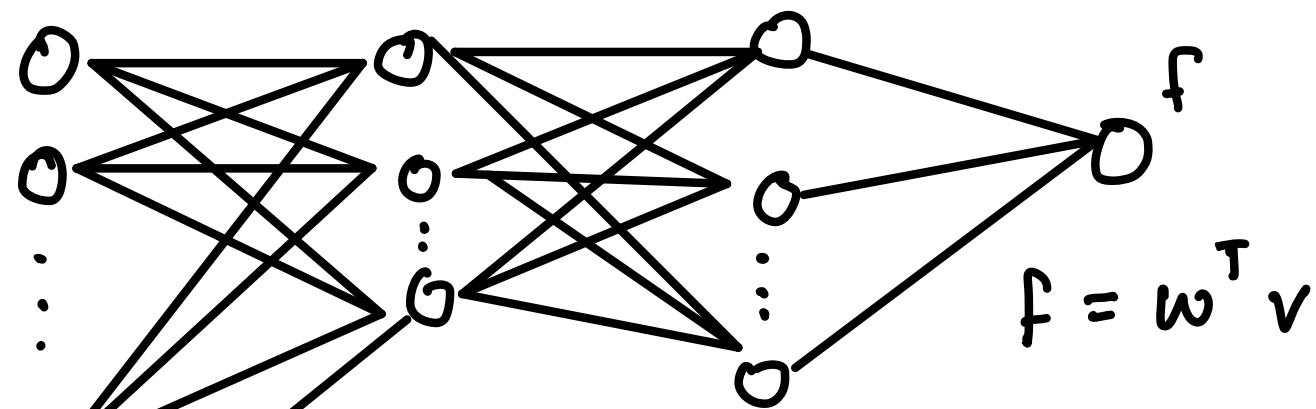
small  $\hat{y}_j \rightarrow$  large loss

large  $\hat{y}_j \rightarrow$  small loss



# Neural Networks

$$x \in \mathbb{R}^d \quad u \in \mathbb{R}^h \quad v \in \mathbb{R}^h$$



$$u = \sigma(W^{(1)} x) \quad v = \sigma(W^{(2)} u)$$

$x \quad h \quad h \times d \quad d \quad h \quad h \times h \quad h$

$$u_i = \langle W_{(i)}^{(1)}, x \rangle + b_i = \sum_{j=1}^d x_j \cdot W_{ij}^{(1)} + b_i$$

$$f = w^T \underbrace{\sigma(W^{(2)})}_{h \times h} \underbrace{\sigma(W^{(1)} x)}_{h \times d}$$

$1 \times h \quad h \times h \quad h \times d$

$$\sigma(v) = \begin{bmatrix} \sigma(v_1) \\ \sigma(v_2) \\ \vdots \end{bmatrix}$$

# Nonlinear Activations!

$$\sigma(z) = z \quad \text{linear}$$

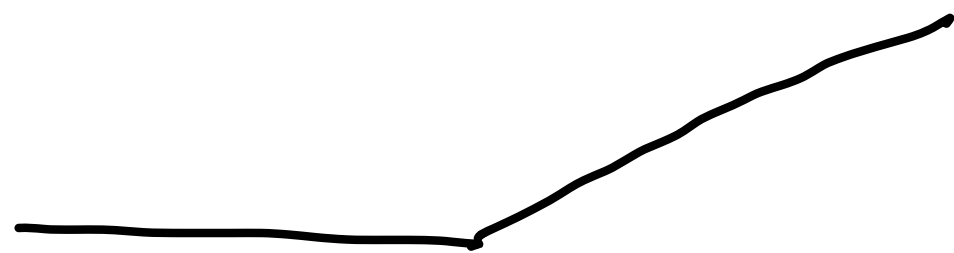
$$\sigma(z) = \frac{\exp(z)}{\sum_{z'} \exp(z')} \quad \text{softmax}$$

$$\sigma(z) = \frac{1}{1 + e^{-z}} \quad \text{sigmoid}$$

$$\sigma(z) = \max(0, z) \quad \text{ReLU}$$

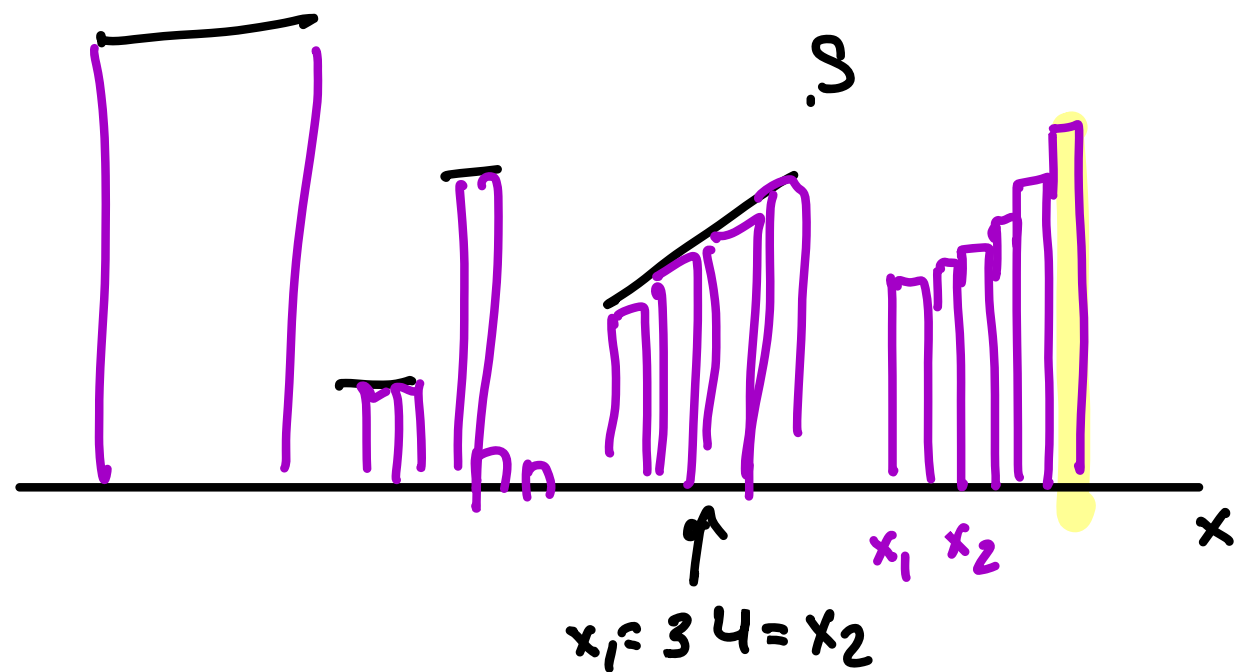
$$\sigma(z) = \text{step}(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ -1 & \text{else} \end{cases}$$

$$\sigma: \mathbb{R} \rightarrow \mathbb{R}$$



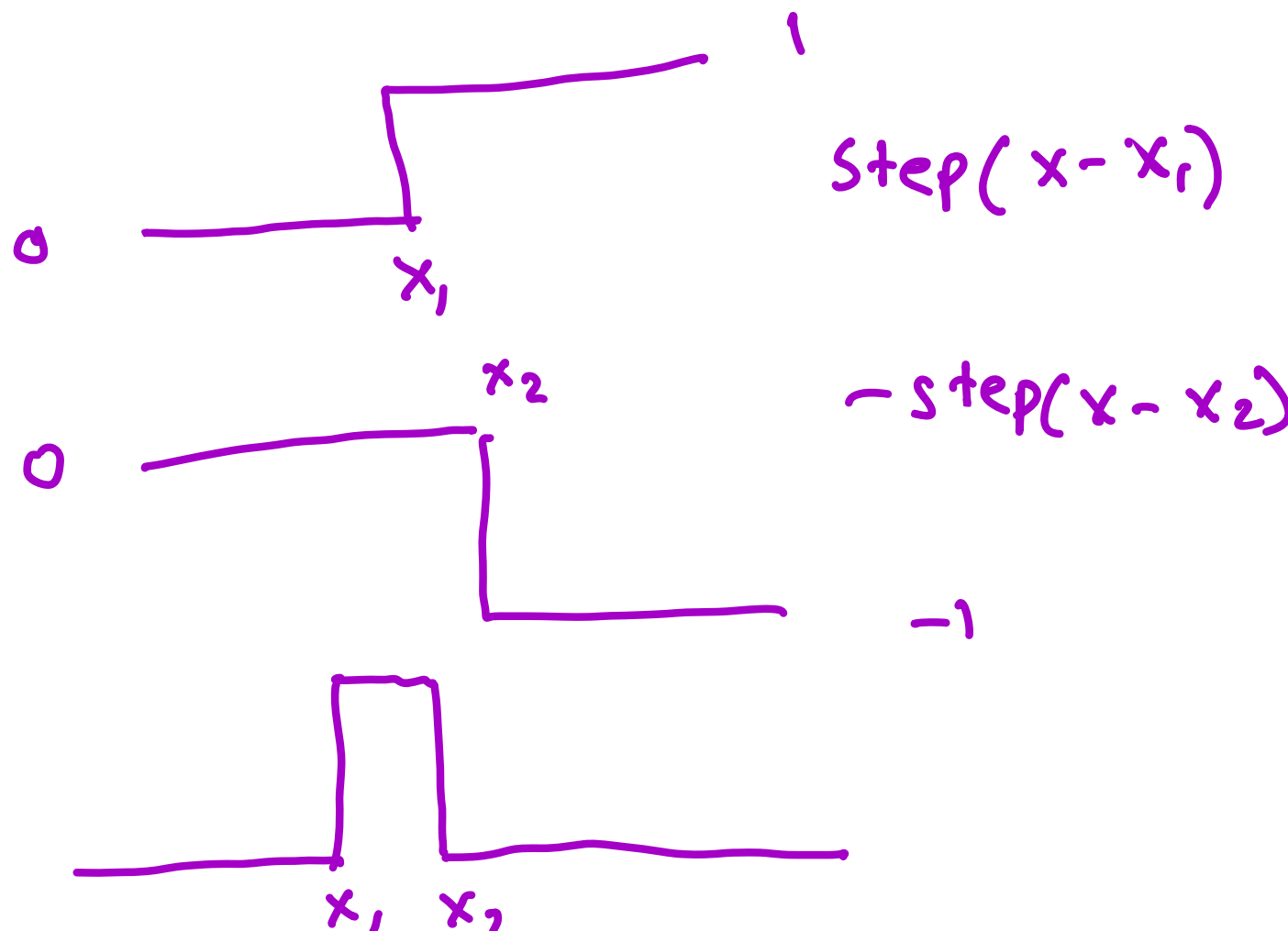
Do neural networks fit everything?

How about this function?

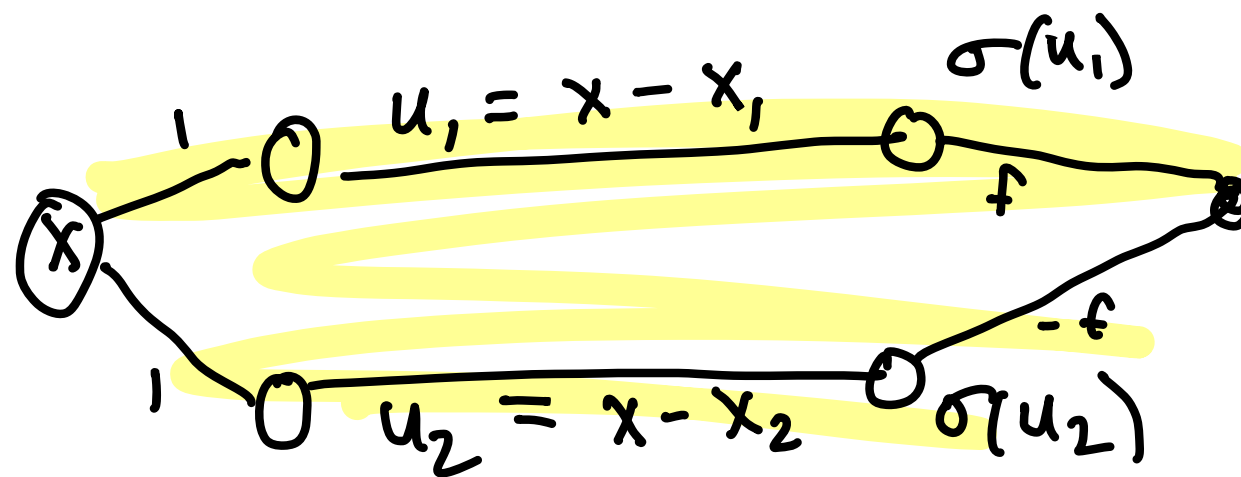


$$\text{step}(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{else} \end{cases}$$

$f(x)$



$$\text{step}(x - x_1) - \text{step}(x - x_2) \quad \sigma = \text{step}$$



$$f(\sigma(u_1) - \sigma(u_2))$$

Good models exist.

How do we find good models?

Gradient Descent!

$\mathcal{L}(w)$

want  $\Delta$  so that

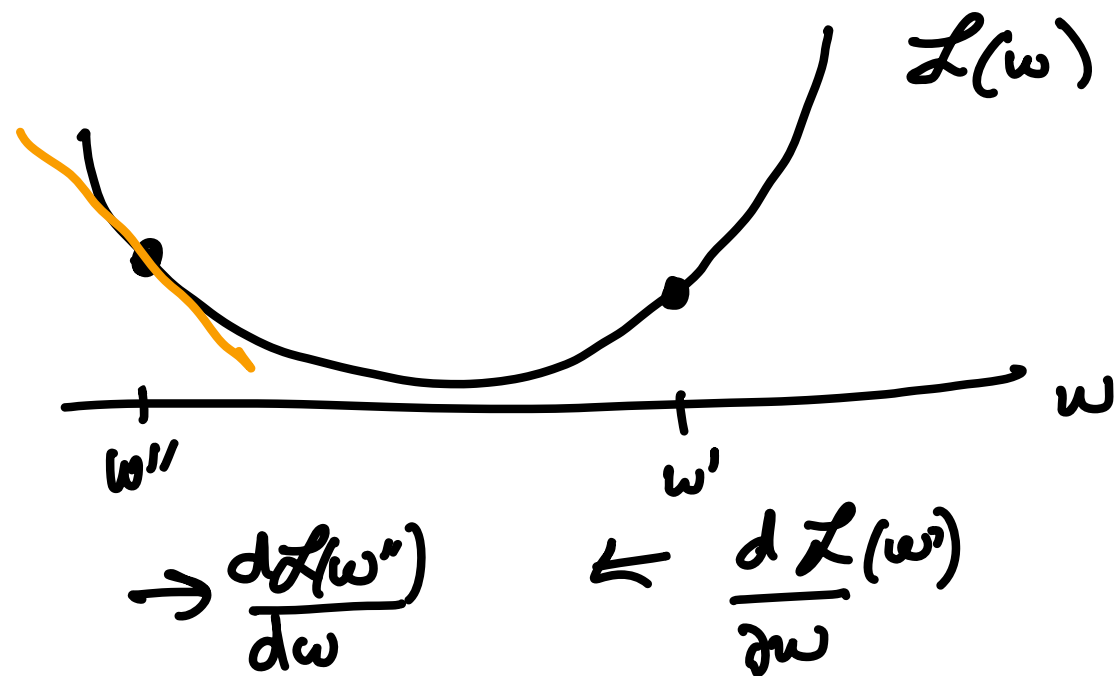
repeat  $\mathcal{L}(w + \Delta) < \mathcal{L}(w)$

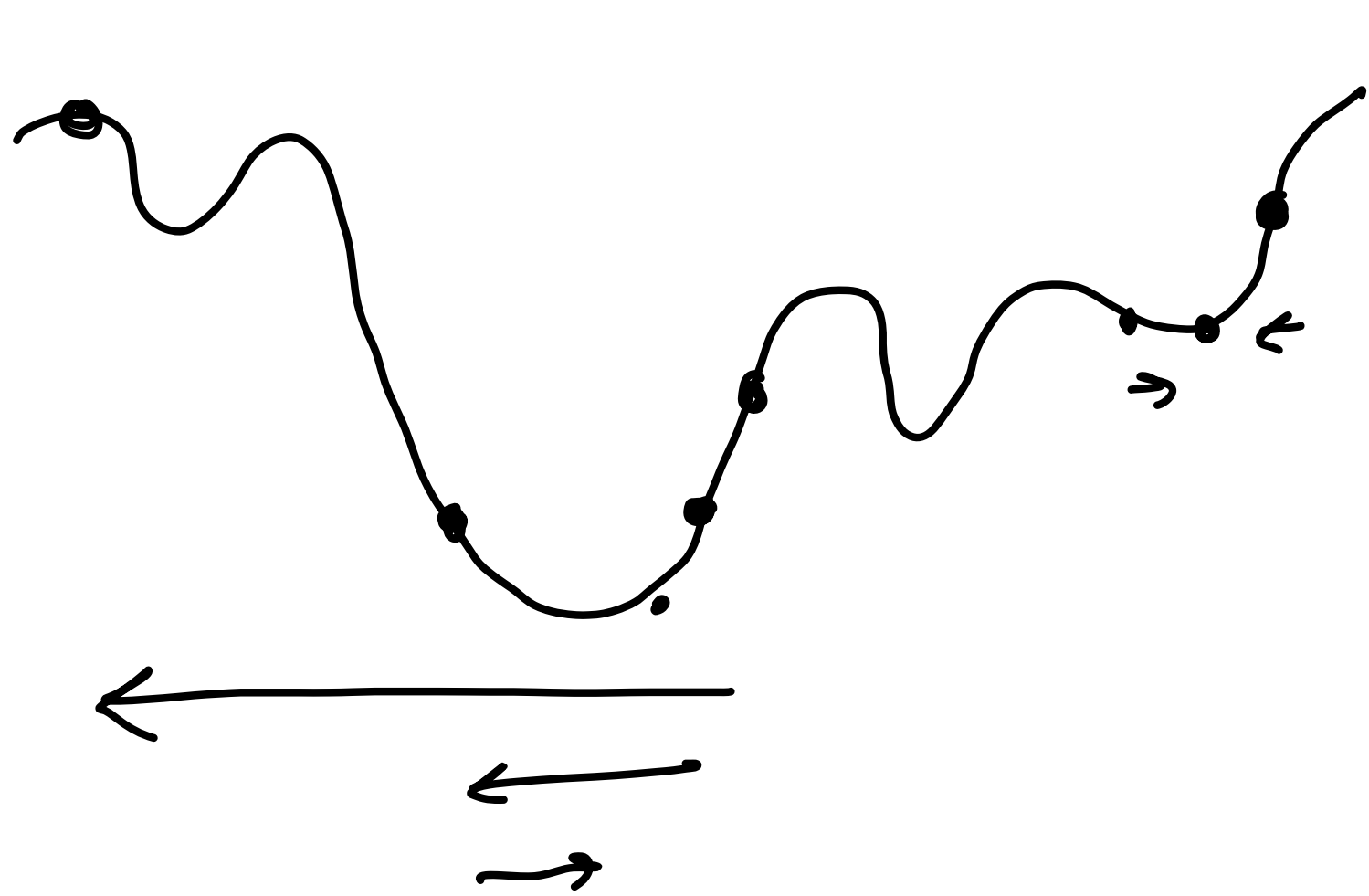
$w \leftarrow w + \Delta$

$$\nabla \mathcal{L}(w) = \begin{bmatrix} \frac{d\mathcal{L}(w)}{dw_1} \\ \vdots \end{bmatrix} \quad w \in \mathbb{R}^d$$

$$w \leftarrow w - \alpha \nabla \mathcal{L}(w)$$

↑  
learning rate





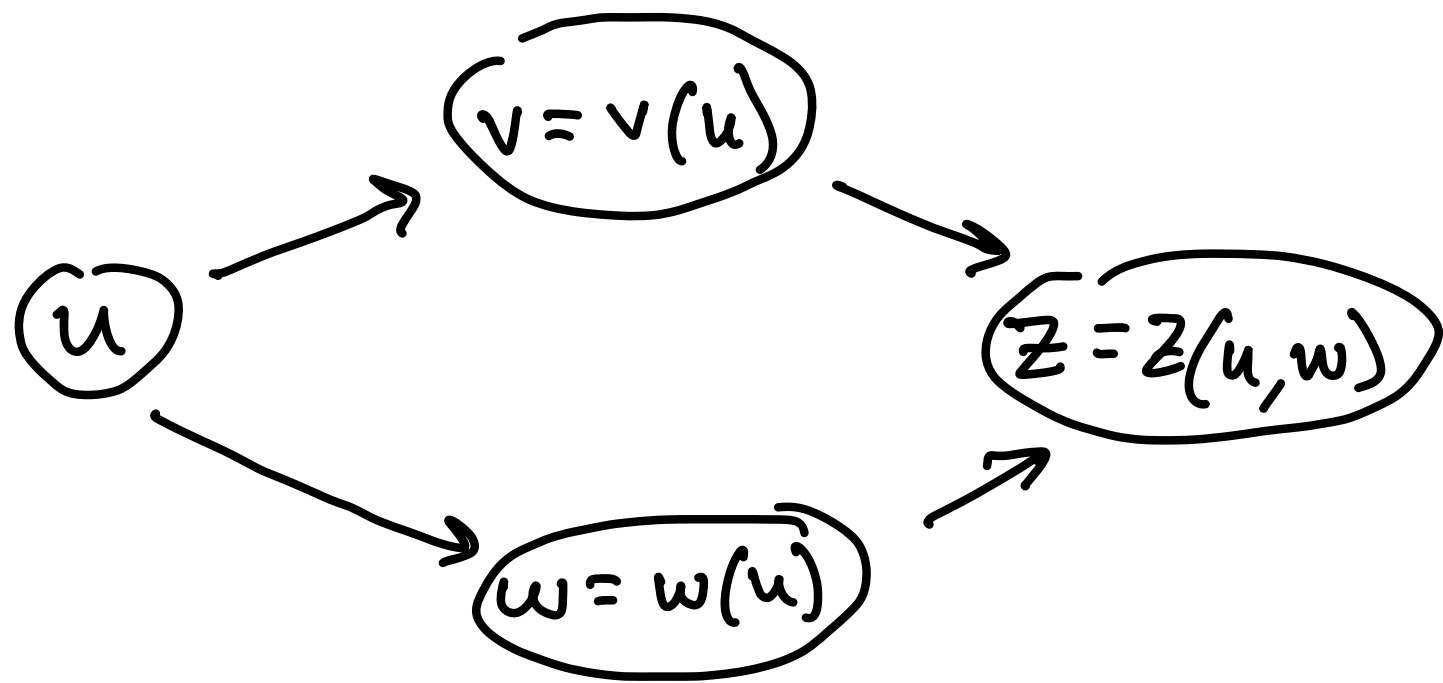
$$z = \omega x + b \quad f(z) = \sigma(z)$$

$$\mathcal{L}(\omega, b) = \frac{1}{2} (y - \sigma(\omega x + b))^2 + \lambda \omega^2$$

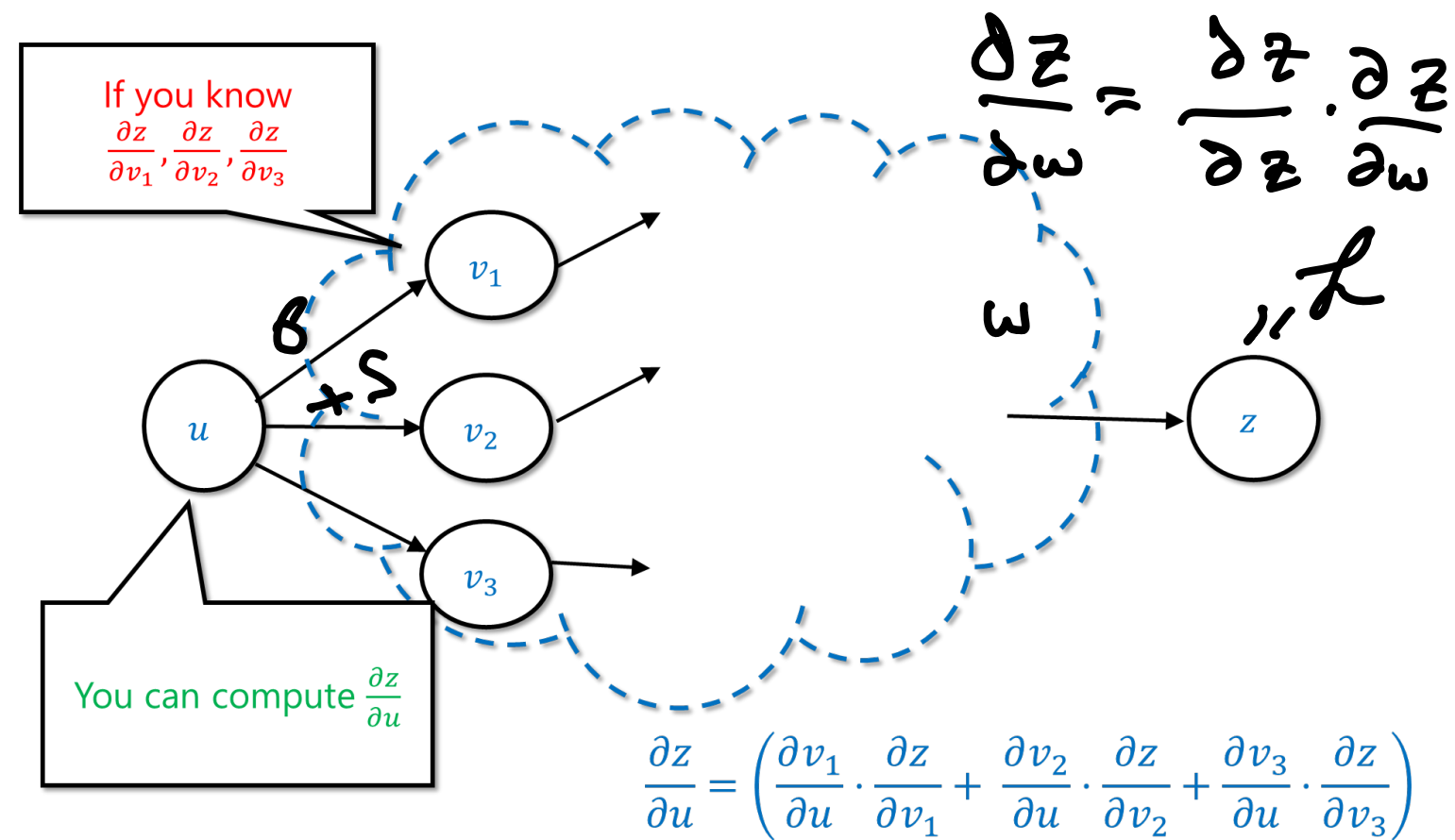
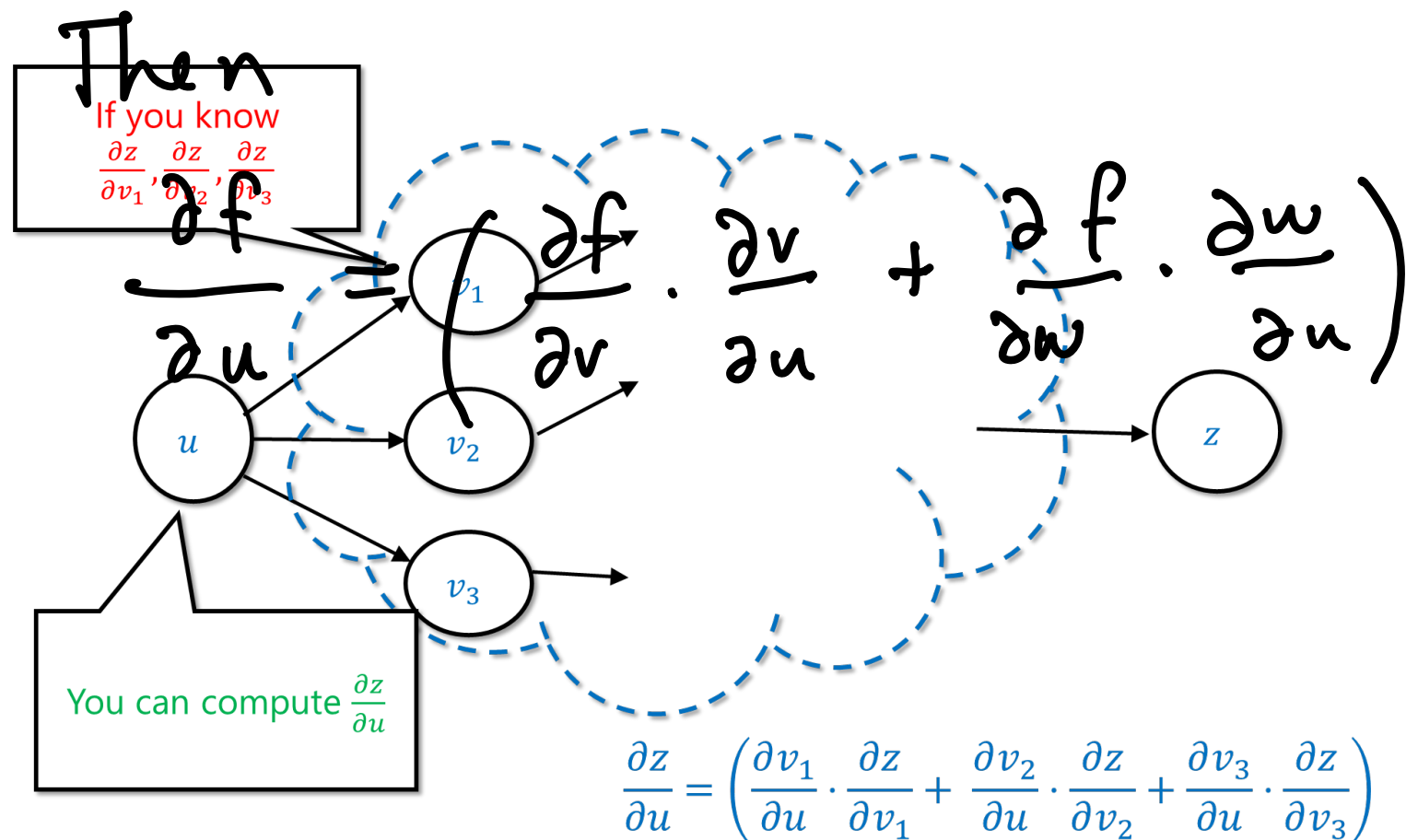
$$\frac{\partial \mathcal{L}}{\partial \omega} = \frac{1}{2} \cdot 2 (y - \sigma(\omega x + b)) \cdot \sigma'(\omega x + b) \cdot x + 2\lambda \omega$$

$$\frac{\partial \mathcal{L}}{\partial b} = \frac{1}{2} \cdot 2 (y - \sigma(\omega x + b)) \cdot \sigma'(\omega x + b) \cdot (-1)$$





$$f(u) = z(v(u), w(u))$$



## Forward pass

for each node:

compute  $v_i$  using  $\text{Parents}(v_i)$

## Backward

for each node  $u$ :

$$\frac{\partial z}{\partial u} = \sum_{v \in \text{children}(u)} \frac{\partial z}{\partial v} \cdot \frac{\partial v}{\partial u}$$

Forward pass

for each node:

compute  $v_i$  using  $\text{Parents}(v_i)$

Backward

for each node  $u$ :

$$\frac{\partial \mathcal{L}}{\partial u} = \sum_{v \in \text{children}(u)} \frac{\partial \mathcal{L}}{\partial v} \cdot \frac{\partial v}{\partial u}$$

$$v = u^2$$

$$\frac{dv}{du} = 2u$$

assuming  $\frac{\partial \mathcal{L}}{\partial v}$  is known

then  $\frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial v} \cdot \frac{\partial v}{\partial u}$