

Exam:

Cumulative

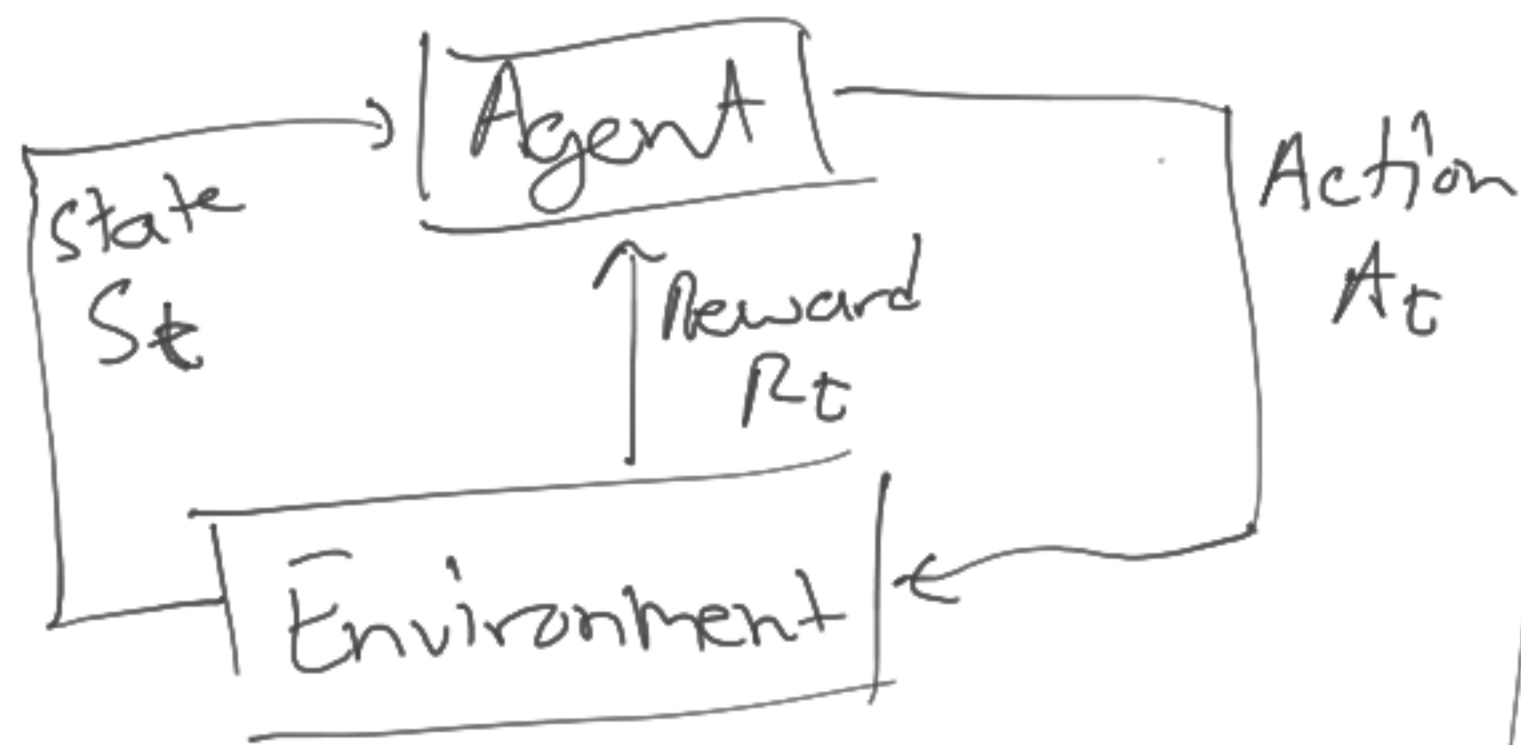
Completable in 2 hour exam slot. (May 11, 2pm-4pm)

I will be available during 2 hour exam slot

Posted 12:01 am May 11, due 11:59 pm May 11.

Approximately same format as midterm.

More focus on RL than other topics.



$t \in \{0, 1, 2, \dots\}$ (time step)

S = set of possible states

A = set of possible actions

S_t = state at time t

A_t = action at time t

R_t = reward at time t

s = some state

a = some action.

$\pi(S, a) = P(A_t = a | S_t = s)$, $\forall t$ ^{remember, R_t generated using π .}

$J(\pi) = \text{how good is } \pi? = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R_t \right]$ ^{semi-colon}

γ = reward discount parameter. $\gamma \in [0, 1]$. ^{expected value. \rightarrow bold E.}

$\pi^* \in \arg \max_{\pi \in \Pi} J(\pi)$ ^{in, not =}

Π^* ^{optimal policy} Π ^{set of all possible policies}

episode = one run ^{starting from $t=0$}

θ = policy parameters (vector).

$\pi_{\theta}(s, a) = P_r(A_t = a | S_t = s; \theta)$ ^{semi-colon.}

$\pi_{\theta}(s, a) = \frac{e^{\sum_j \theta_j^a \phi_j(s)}}{\sum_{a'} e^{\sum_j \theta_j^{a'} \phi_j(s)}}$ ^{softmax selection: action}

$\theta = [\theta^1, \theta^2, \dots, \theta^{|A|}]$

$\theta^a = [\theta_1^a, \theta_2^a, \dots, \theta_{\text{numFeatures}}^a]$

$V^{\pi}(s) = \text{value of state } s \text{ under policy } \pi$

$= \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k} \mid \pi \right]$

w = weights/parameters for value function estimate

v_w = estimate of v^{π} using weights w .

$v^{\pi}(s) \approx v_w(s) \quad \forall s \in S$

$$\delta_t = R_t + \gamma v_w^{\pi}(S_{t+1}) - v_w^{\pi}(S_t) \quad \text{OR } \square$$

for each episode

for each time t

Agent observes S_t

Agent selects action A_t using π_θ

Environment responds by transitioning from S_t to S_{t+1} and emitting reward R_t

$\delta_t = R_t + \gamma v^\pi(S_{t+1}) - v^\pi(S_t)$ → critic judging the actor.
 $\theta_i, \phi_i \leftarrow \theta_i + \alpha \delta_t \frac{\partial \ln(\pi_\theta(S_t, A_t))}{\partial \theta_i}$ → actor update

$w_j, \omega_j \leftarrow w_j + \beta \delta_t \frac{\partial v_w(S_t)}{\partial w_j}$ → critic update

$\sum_{k=0}^{\infty} \gamma^k R_{t+k}$ → R_t and S_{t+1} are all we have.

Run the episode

convert Eng → math

get rid of "if $\sum R_t$ is big"

A_t only responsible for rewards that happen after A_t .

move update into episode using TD error δ_t .

learn v^π
 Move v^π update into episode using δ_t .

place here

if $\sum_{t=0}^{\infty} \gamma^t R_t$ is big

for each time t

Make action A_t more likely in state S_t

$\theta_i, \phi_i \leftarrow \theta_i + \alpha \frac{\partial \ln(\pi_\theta(S_t, A_t))}{\partial \theta_i}$ → How to change θ_i to increase $\pi_\theta(S_t, A_t)$

$\theta_i, \phi_i \leftarrow \theta_i + \alpha \left(\sum_{k=0}^{\infty} \gamma^k R_{t+k} \right) \frac{\partial \ln(\pi_\theta(S_t, A_t))}{\partial \theta_i}$

optional. (comm)

Actor-Critic

cut all of this.

if $\sum_{t=0}^{\infty} \gamma^t R_t$ is small

for each time t

Make action A_t less likely in state S_t

$\theta_i, \phi_i \leftarrow \theta_i - \alpha \frac{\partial \ln(\pi_\theta(S_t, A_t))}{\partial \theta_i}$

replace w/ $R_t + \gamma v^\pi(S_{t+1}) - v^\pi(S_t)$
 what happened after step t what we expected at time t .

learn from the episode

| input | output |
|----------|--|
| S_0 | $\sum_{k=0}^{\infty} \gamma^k R_{0+k}$ |
| S_1 | $\sum_{k=0}^{\infty} \gamma^k R_{1+k}$ |
| \vdots | \vdots |
| | $R_t + \gamma v_w(S_{t+1})$ |

For each time t :
 $\theta_j, w_j \leftarrow w_j + \beta \left(\sum_{k=0}^{\infty} \gamma^k R_{t+k} - v_w(S_t) \right) \frac{\partial v_w(S_t)}{\partial w_j}$

Supervised Learning

x_i = input for i th point
 y_i = output/label for i th point
 n = number of points
 \hat{y}_i = agent's prediction of y_i from x_i .

non-parametric: NN
 kNN
 WKNN

parametric:

w_k = model parameters after k updates.

$$f_{w_k}(x_i) = \hat{y}_i$$

Linear parametric model:

basis function ϕ_j

$\phi_j(x_i)$ = j th feature for input x_i .

$$f_{w_k}(x_i) = \sum_j w_{kj} \phi_j(x_i)$$

$l(w_k)$ = loss function.
 (smaller = better w_k)

least squares loss:

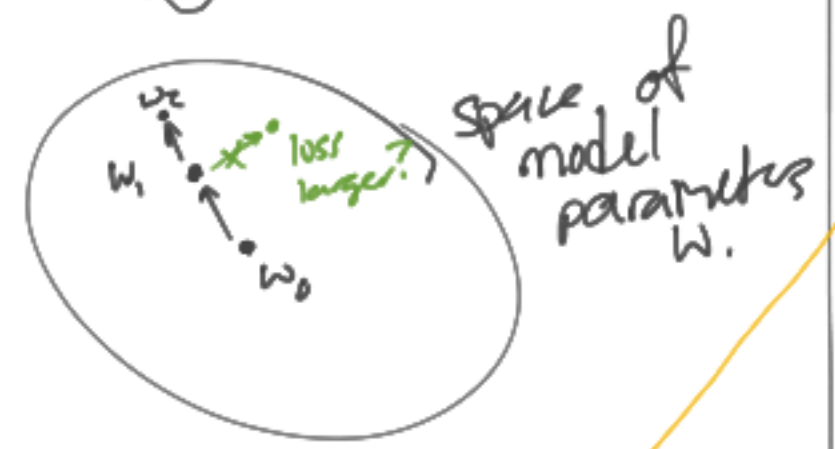
$$l(w_k) = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Other forms:

$$l(w_k) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$l(w_k) = \frac{1}{2n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Hill climbing:
(descent)



~~HA~~ $\sin(x)$ is Lipschitz, $L=1$
~~HA~~ $\sin(x^2)$ is not Lipschitz.

Gradient descent:
 $w_{k+1,j} = w_{k,j} - \alpha \partial l(w_k) / \partial w_{k,j}$

α = step size.

α_k = step size for k th update.

Remember all hyperparams, like when to stop!

Convergence:

- Square summable, not summable

$$\sum_{k=0}^{\infty} \alpha_k^2 < \infty$$

$$\sum_{k=0}^{\infty} \alpha_k = \infty$$

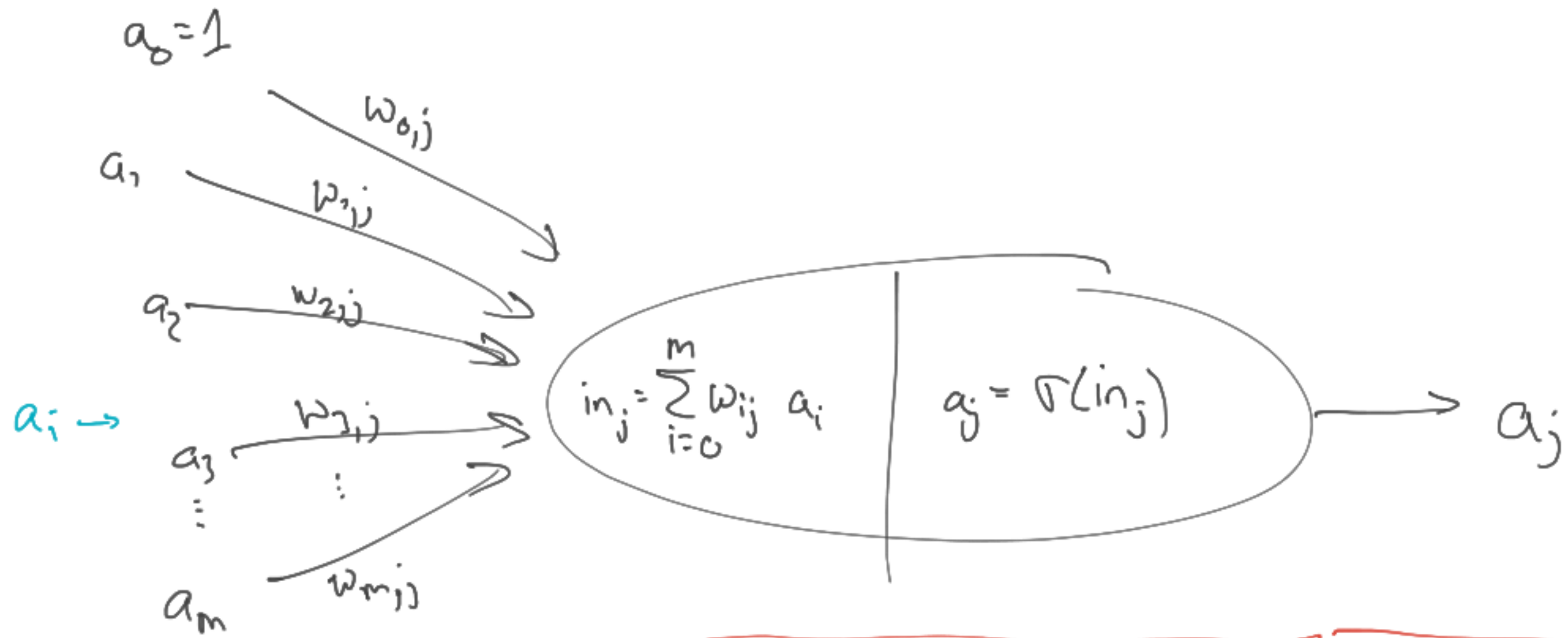
Lipschitz continuity

= If differentiable function, max slope is L .

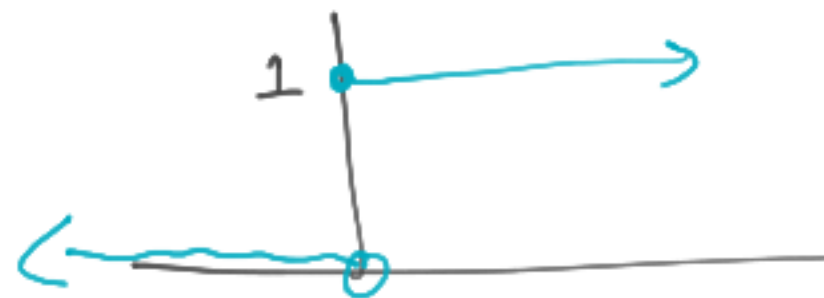
Don't forget to normalize inputs!

Each feature \rightarrow mean = 0

\rightarrow variance = 1



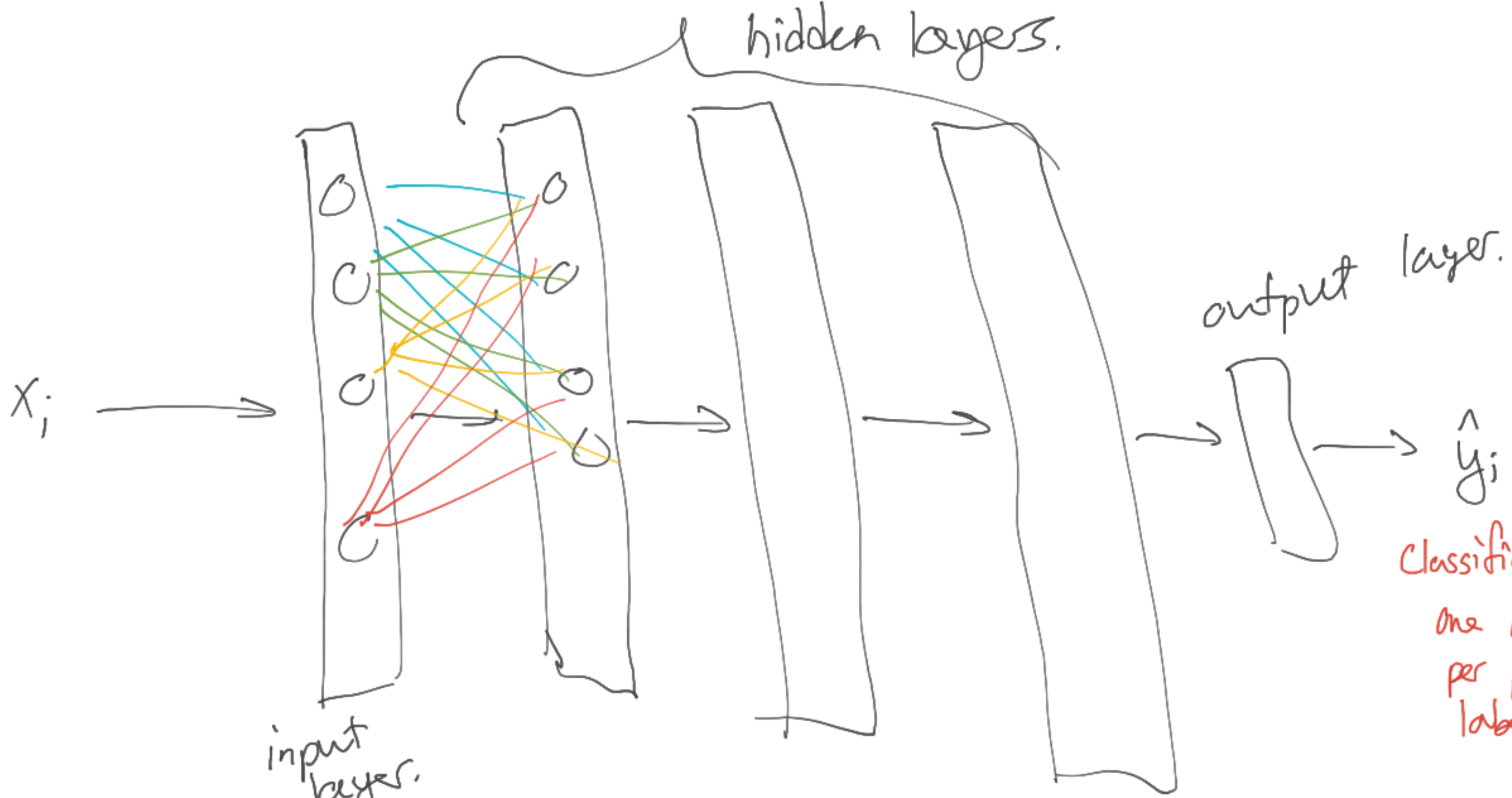
$$\sigma(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{otherwise} \end{cases}$$



$\sigma(z)$ = sigmoid
 \rightarrow logistic function

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

$$\frac{\partial \sigma(z)}{\partial z} = \sigma(z)(1 - \sigma(z))$$



Fully connected feedforward artificial neural network

Classification:
one output
per possible
label.

Backpropagation:

Old notation

$$\frac{d l(w_k)}{d w_{kj}} =$$

(

$$\frac{\partial f_{w_k}(x_i)}{\partial w_{kj}})$$

New notation

$$\frac{\partial f_w(x)}{\partial w_{ij}} =$$

(chain rule over, and over)

$i = \text{prev layer}$

$j = \text{cur layer}$

$k = \text{next layer}$

Other:

Vanishing gradients

over-fitting

train, validation, test

adaptive step sizes

generalization bounds (safety)

- Hoeffding's inequality
applied to $l(w_k)$

- using test data!

fairness

psych/neuro

ethics

philosophy of mind