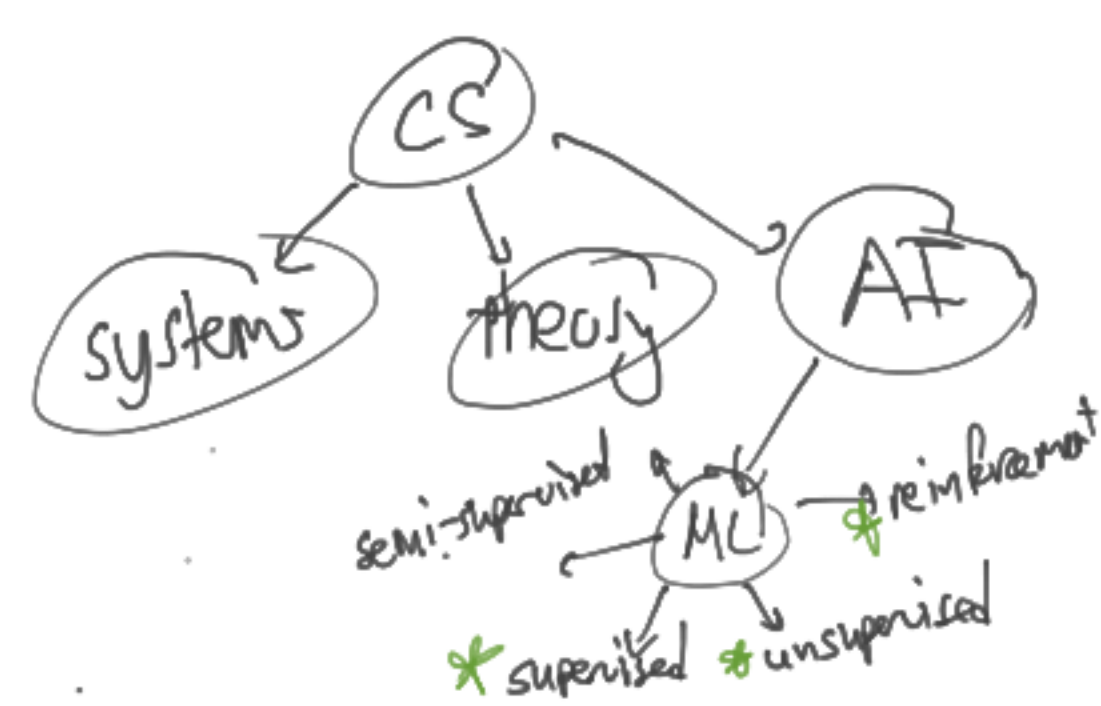
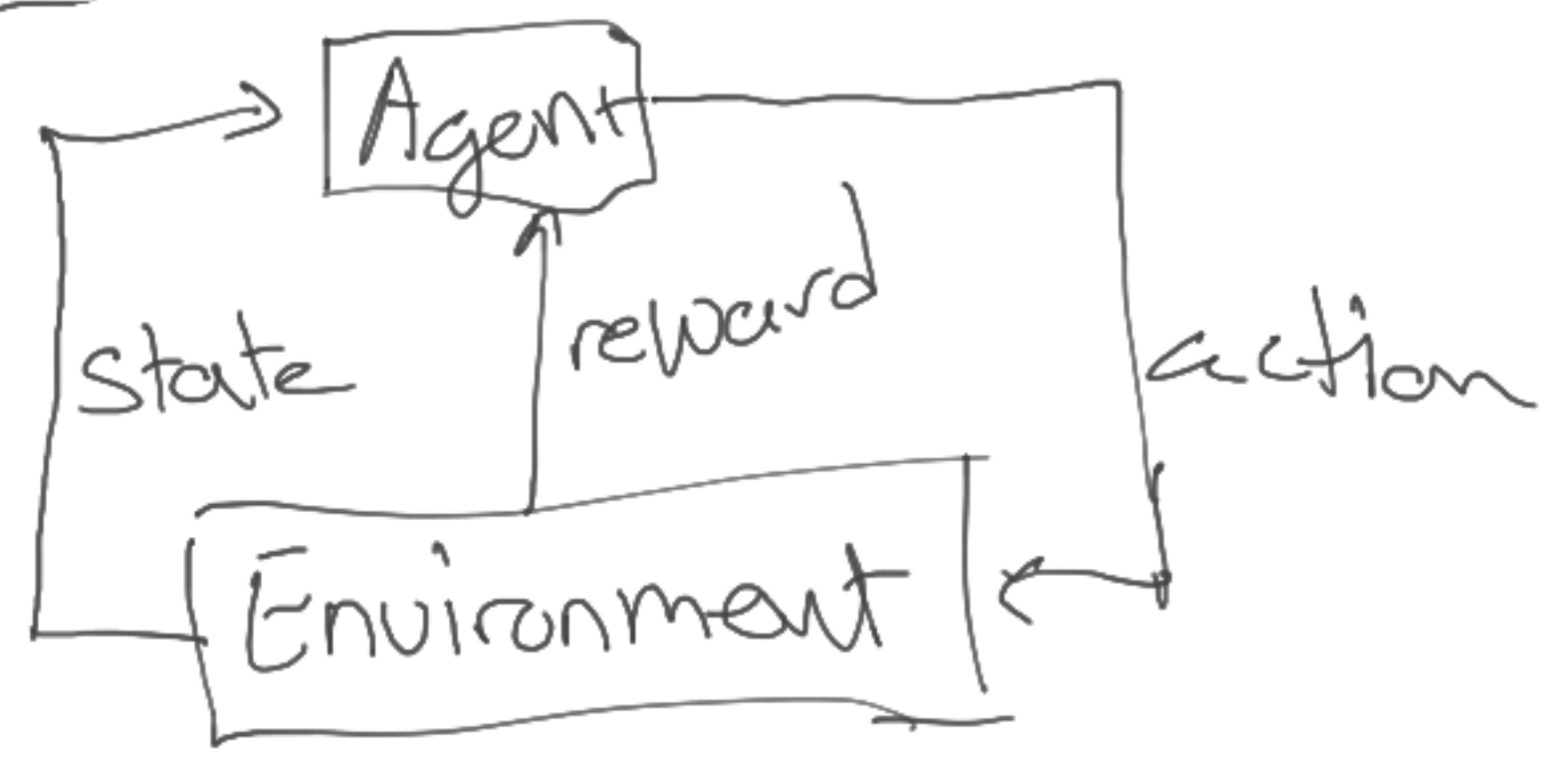


What is reinforcement learning (RL)?

"Reinforcement learning is an area of machine learning, inspired by behaviorist psychology, concerned with how an agent can learn from interactions with an environment."

- Wikipedia / Sutton & Barto / Phil
1998

Agent - Environment Diagram



Neuroscience:

How do animals learn?

A specific agent

or set of agents.

- The study of some examples of learning and intelligence!

RL

(ML)

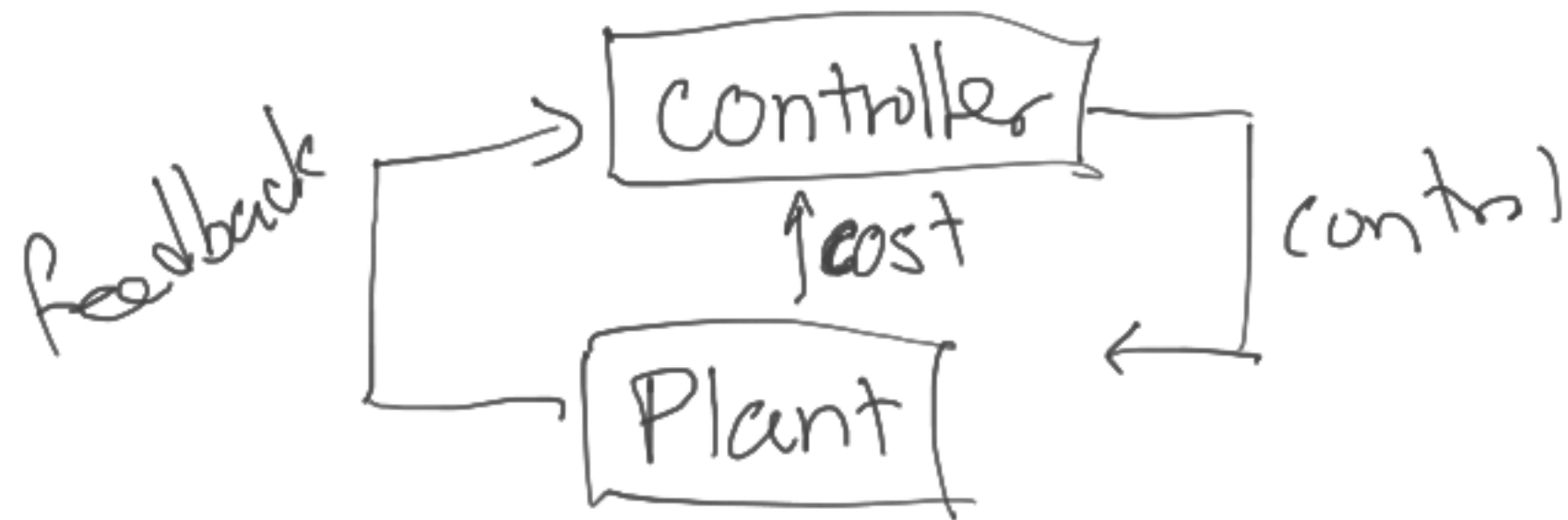
: How can we make an agent that learns?

- The study of learning & intelligence. (in general)

Dopamine \approx Temporal difference error.

Two most related fields

- Operations research
- Control (adaptive / classical)



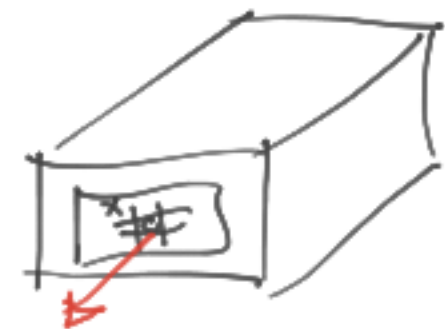
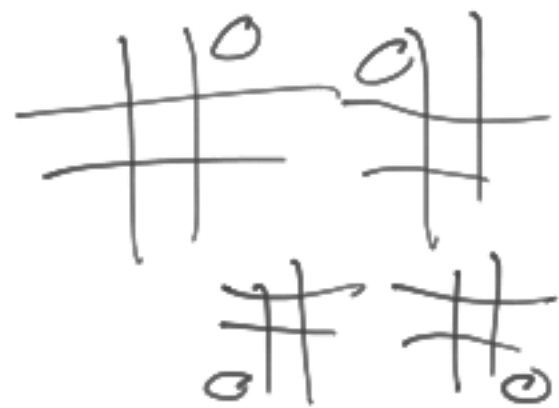
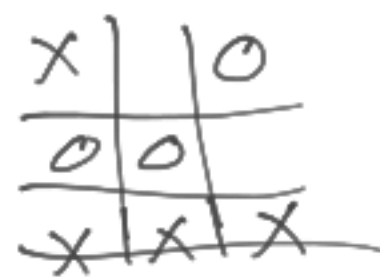
- Main difference is that these fields typically assume the environment (Plant) can and should be directly approximated.



Key properties:

- 1) Evaluative feedback
↳ "This is how good the outcome was," (not instructional)
not "this is what you should have done."
- 2) Sequential.
↳ ~~no~~ no data set!

MENACE (Donald Michie, 1961)



304 matchboxes.

4 bead colors.

↳ possible moves.

Win: Put back bead, +3 more of same color.

Loss: Remove beads.

Tie: Put back bead, +1 more of same color.
(Draw)

Take:

- Act at least slightly random.
- Good outcome
↳ take chosen actions more often
- Bad outcome
↳ take chosen actions less often.

Act randomly ???

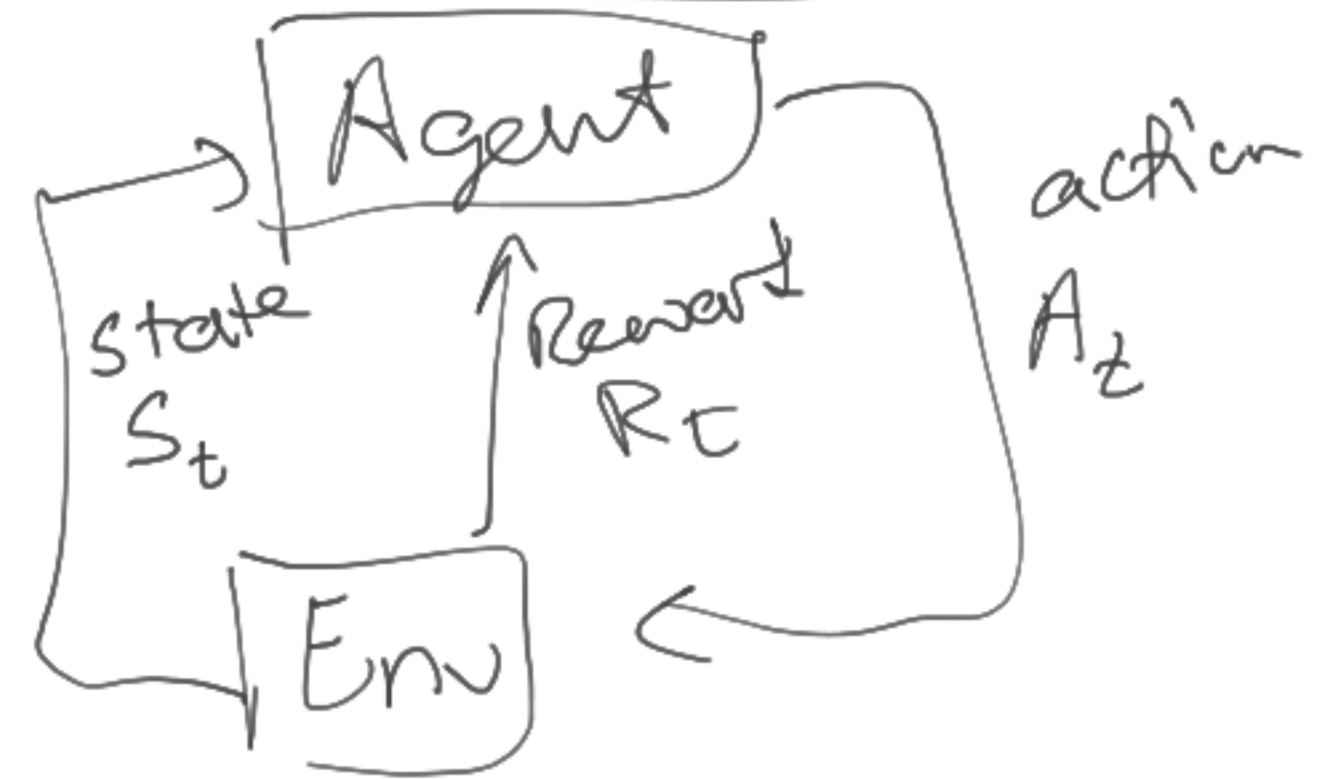
really exploration is necessary, and randomness is one way to explore

EMPSCI
687

- "Necessary" for learning
- "exploration" vs "exploitation"

↳ taking actions that think are suboptimal

↳ choosing actions that think are best.



t : time step.
 $t \in \{0, 1, 2, \dots\}$

S_t : State at time t .
★ (Agent's observation of state)

A_t : Action at time t .

R_t : Reward at time t .

$R_t \in \mathbb{R}$

A policy π is a way of selecting actions.

$$\pi(s, a) = \Pr(A_t = a \mid S_t = s)$$

~~$A_t = \pi(S_t)$~~

- many policies - some good some bad.

Agent's goal: find a policy that maximizes the expected amount of reward the agent receives.
 ↳ not deterministic!

Objective function: $J(\pi) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R_t; \pi \right]$

~ Could be infinite! $\rightarrow J(\pi)$

- Doesn't discount based on time
 ↳ care about short term

Means "assuming R_t generated using π "

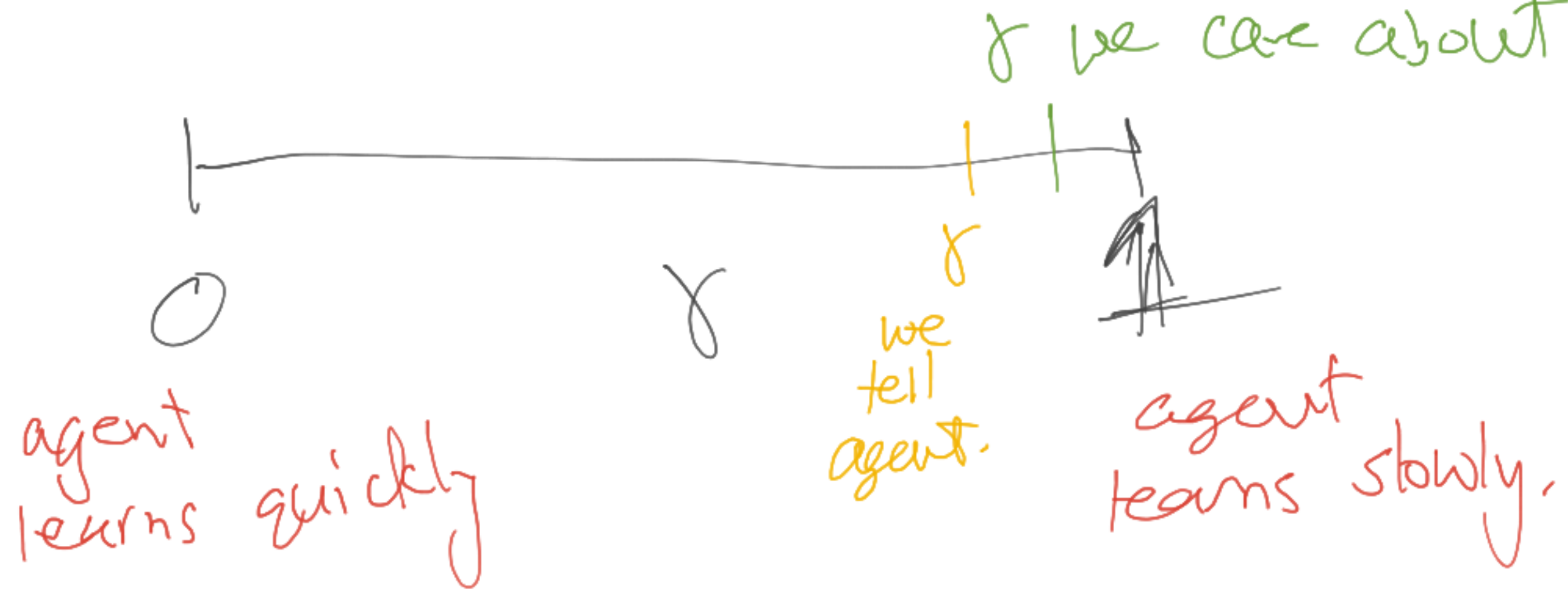
YATP
 ↳ yet another hyperparameters.

$\gamma \in [0, 1]$
 ↳ care about long term.

$$\gamma^0 R_0 + \gamma^1 R_1 + \gamma^2 R_2 + \dots$$

$$1 R_0 + .5 R_1 + .25 R_2 + \dots$$

Typically $\gamma \approx .9 \leftrightarrow 1$



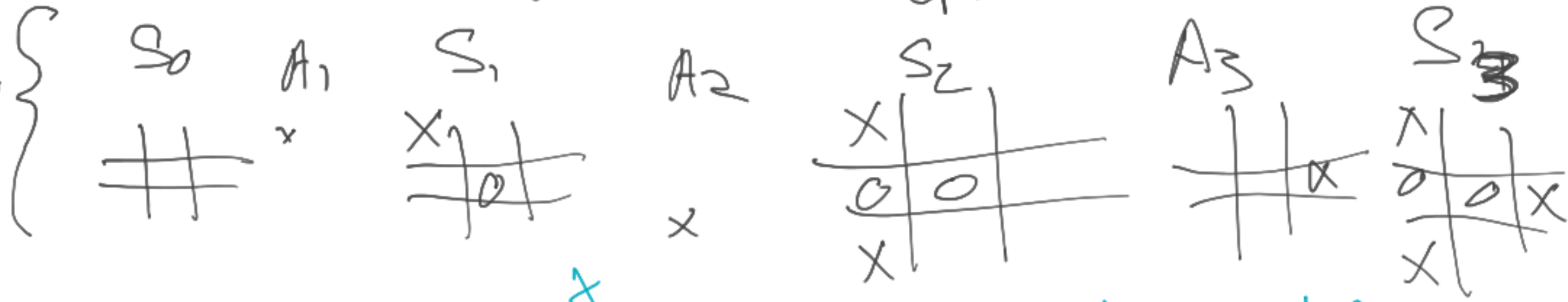
Agent's goal is to find or approximate

$$\pi^* \in \arg \max_{\pi} J(\pi)$$

π^* is called an optimal policy.

"Independent" trials: each starting from $t=0$,
 are called "episodes."

one episode.



next episode.

The dynamics of env do not change across trials, e.g. actions in one episode do not impact environment in later episodes.

Way agent acts (agent's policy) can change across episodes.

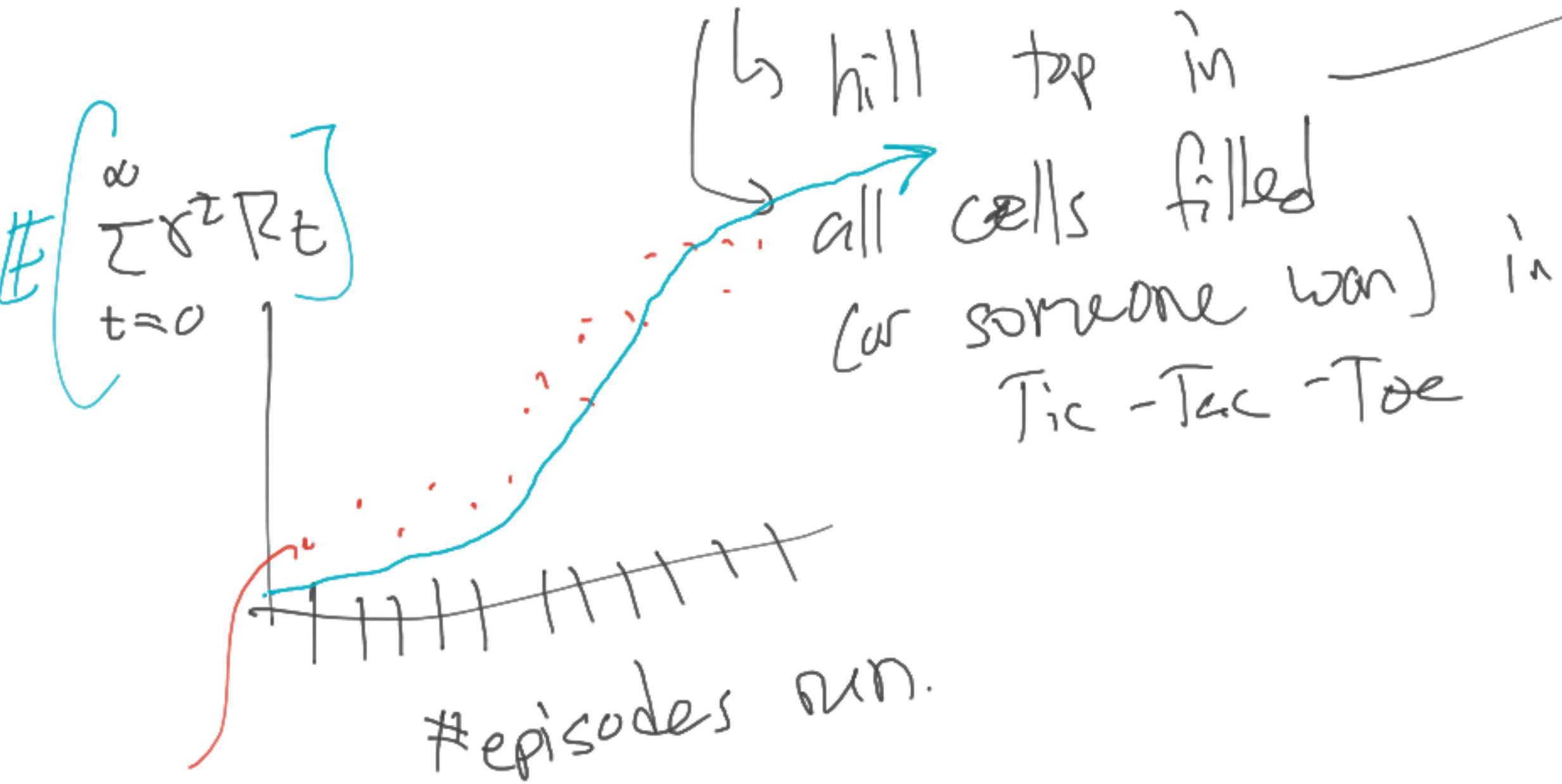
$(S_t, A_t, R_t)_{t=0}^{\infty}$ independent rand variables across episodes.

Episodes end when the environment reaches a state called a "terminal state"

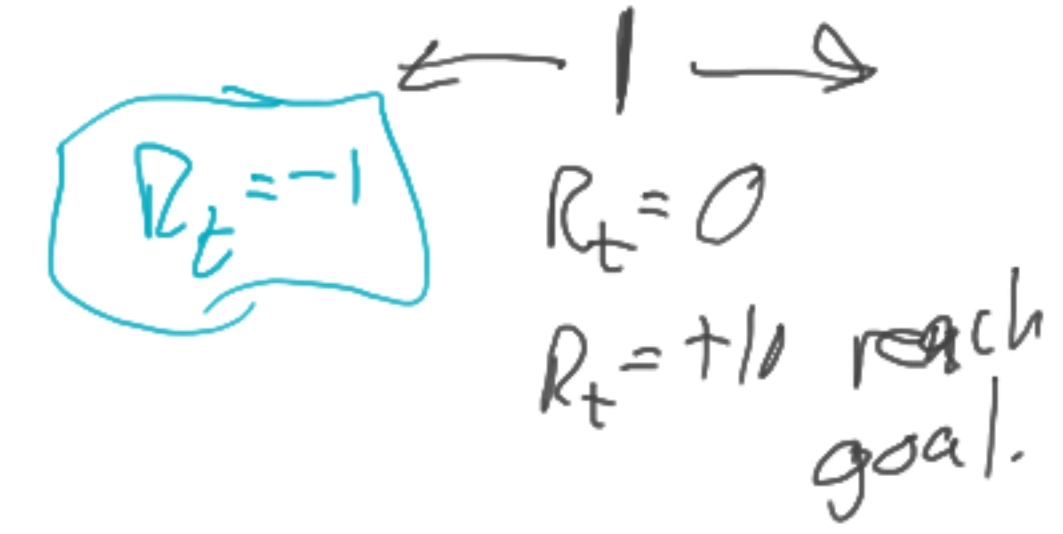
Mountain Car



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



one trial. "lifetime"



Flow:
 $(S_0, A_0, R_0, S_1, A_1, R_1, \dots)$

→ For all times after a terminal state, $R_t = 0$.

Grid world:

$A_t \in \{\text{up, down, left, right}\}$

$S_t = (2, 4)$

$S_t = \emptyset$

$R_t =$

$t=10$
 $t=20$

$S_t = (2, 4)$

$A_t = \text{right}$

$S_{t+1} = (3, 4)$

$R_t = 4$

start S_0

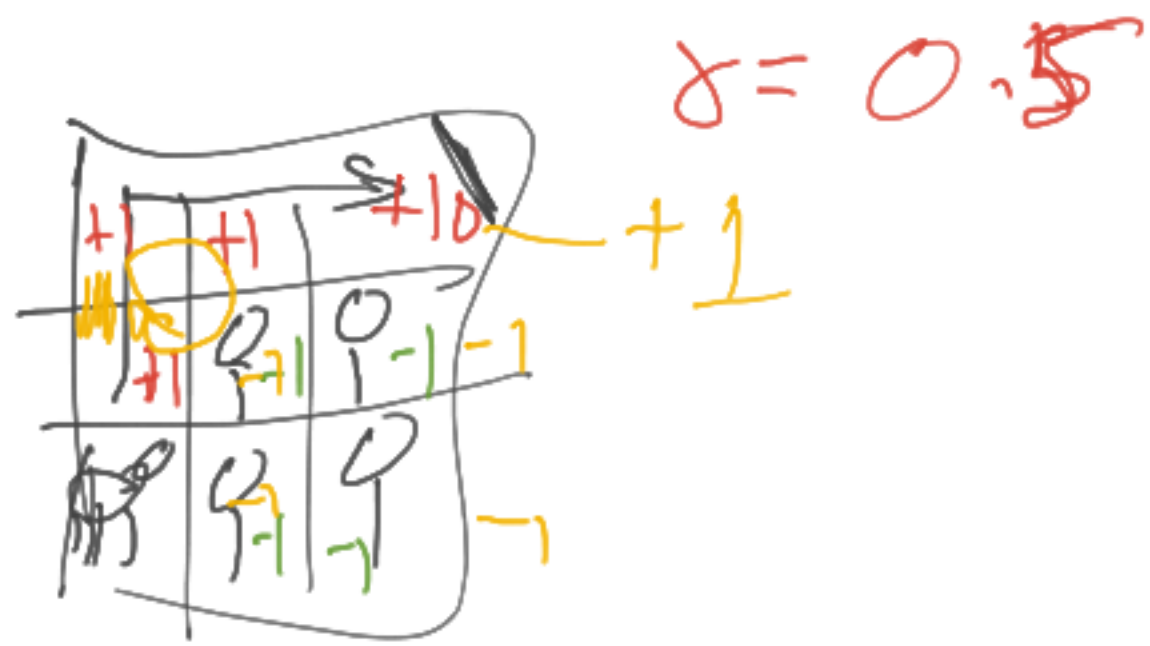


reward for entering a state.

terminal state.

R_t design flower bed

How to define R_t is yet another hyperparameter.



→ Give rewards for the outcomes you want, not for how you think the agent should achieve the outcome.

"reward shaping" Andrew Ng.

How represent π ?

- Simplifying assumption: A_t has a small finite number of possible values.

$S_t \rightarrow \boxed{\cdot} \rightarrow \boxed{\cdot} \rightarrow \boxed{\cdot} \rightarrow \boxed{\cdot} \rightarrow \left. \begin{matrix} \rightarrow o_1 \\ \rightarrow o_2 \\ \vdots \end{matrix} \right\}$ one output per possible action.

Softmax to select the action:

$$\pi(s, a) = \Pr(A_t = a | S_t = s) = \frac{e^{o_a}}{\sum_i e^{o_i}}$$

Supervised learning

RL

model



policy

w

"model parameters"



θ

"policy parameters"

f_w



π_θ

deterministic
for regression.

↳ stochastic, like
classification
(uses softmax)

ℓ



$-J$

$$l(w_k) = \frac{1}{n} \sum_{i=1}^n (y_i - f_{w_k}(x_i))^2$$

$$\frac{\partial l(w_k)}{\partial f_{w_k}(x_i)} = \frac{1}{n} \sum_{i=1}^n 2(y_i - f_{w_k}(x_i)) \frac{\partial}{\partial f_{w_k}(x_i)} (y_i - f_{w_k}(x_i))$$

$$\frac{\partial l(w_k)}{\partial w_{kj}}$$

$$\frac{\partial l(w_k)}{\partial f_{w_k}(x_i)}$$

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

$$= \frac{1}{n} \sum_{i=1}^n 2(y_i - f_{w_k}(x_i)) x_{ij}$$

1c

$$f_{w_k}(x_i) = \sum_{\beta} w_{k\beta} x_{i\beta}$$

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

$$\frac{\partial f_{w_k}(x_i)}{\partial w_{kj}} = \frac{\partial}{\partial w_{kj}} \sum_{\beta} w_{k\beta} x_{i\beta}$$

$$= \sum_{\beta} \frac{\partial}{\partial w_{kj}} w_{k\beta} x_{i\beta} \quad \begin{matrix} j=2 \\ \beta=1 \end{matrix}$$

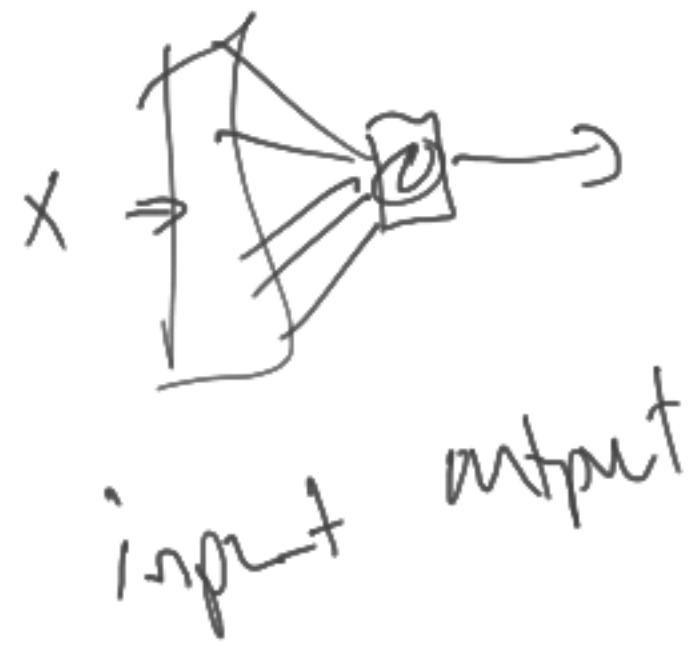
$j=2$

$$= \cancel{\frac{\partial}{\partial w_{k1}} w_{k1} x_{i1}} + \frac{\partial}{\partial w_{k2}} w_{k2} x_{i2} + \cancel{\frac{\partial}{\partial w_{k3}} w_{k3} x_{i3} \dots}$$

$= 0$ $= x_{i2}$ $= 0$

$$= x_{i,2}$$

$$f_{w_k}(x_i) = \sum_{\beta} w_{k\beta} x_{i\beta}$$



linear model
weight, divided
by 7.4.

$$in_j = \sum_i w_{ij} x_i$$

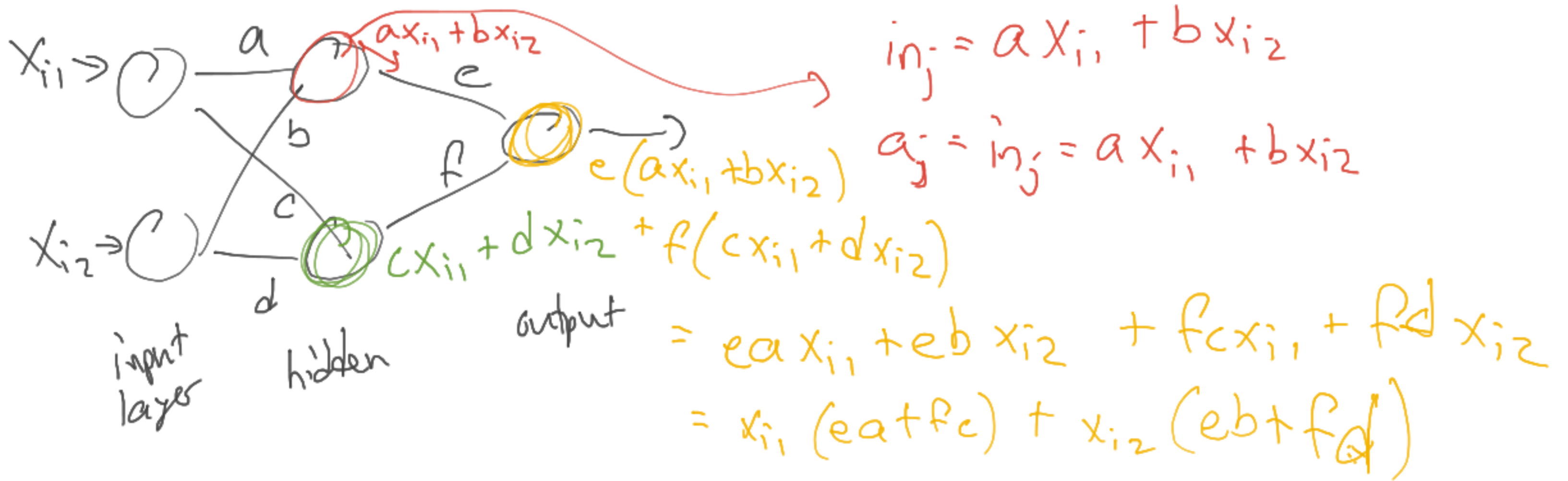
$$a_j = \sigma(in_j)$$

$$a_j = 7.4 in_j \leq 0.5$$

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

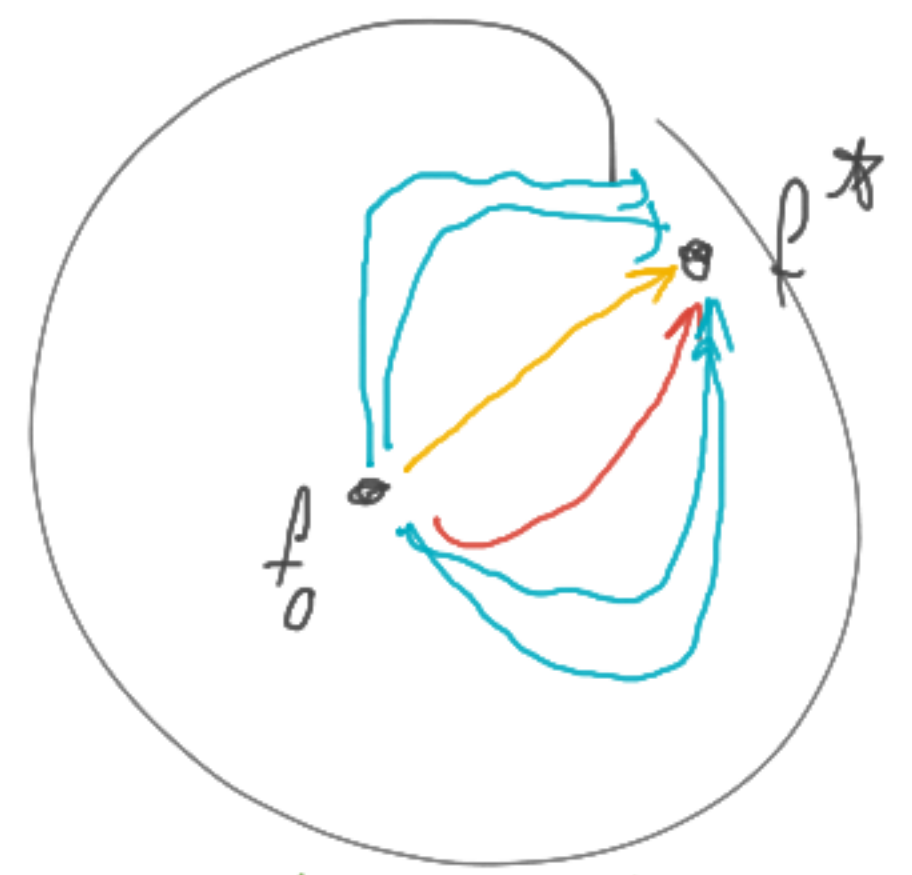
≤ 0.5

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



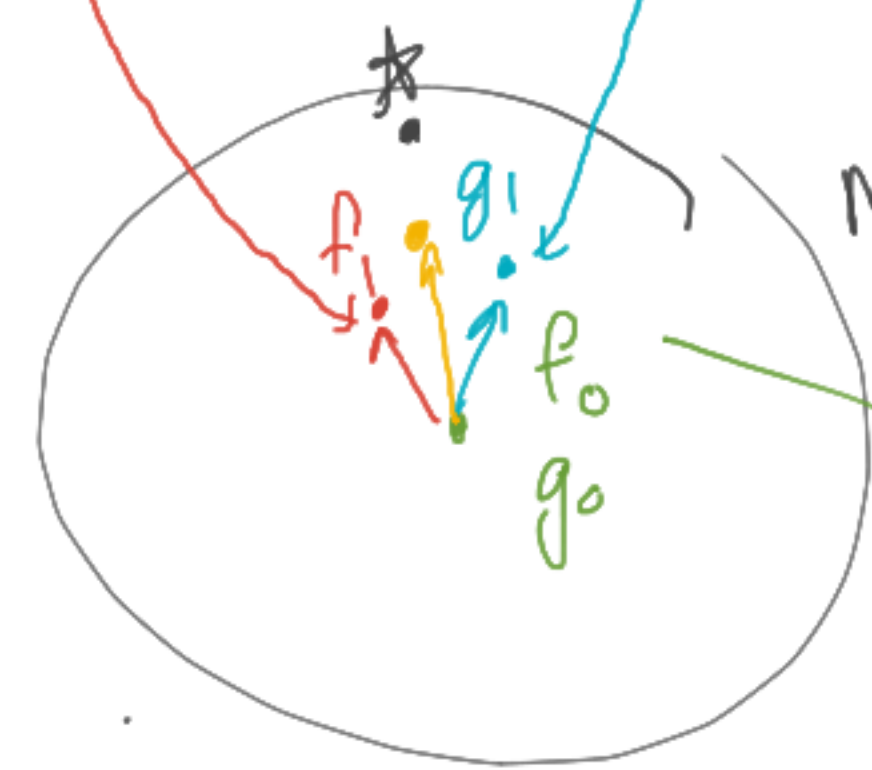
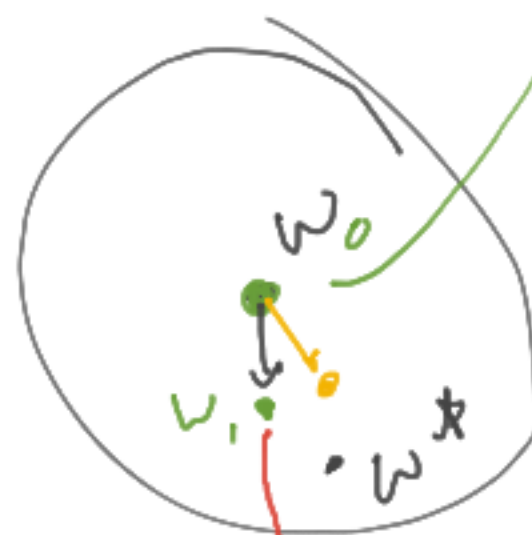
$$f_{w_k}(x_i) = \underbrace{w_{k1}}_{ea+fc} x_{i1} + \underbrace{w_{k2}}_{eb+fd} x_{i2}$$

$$w_0 = \begin{pmatrix} 0 & 0 \\ w_{01} & w_{02} \end{pmatrix}$$

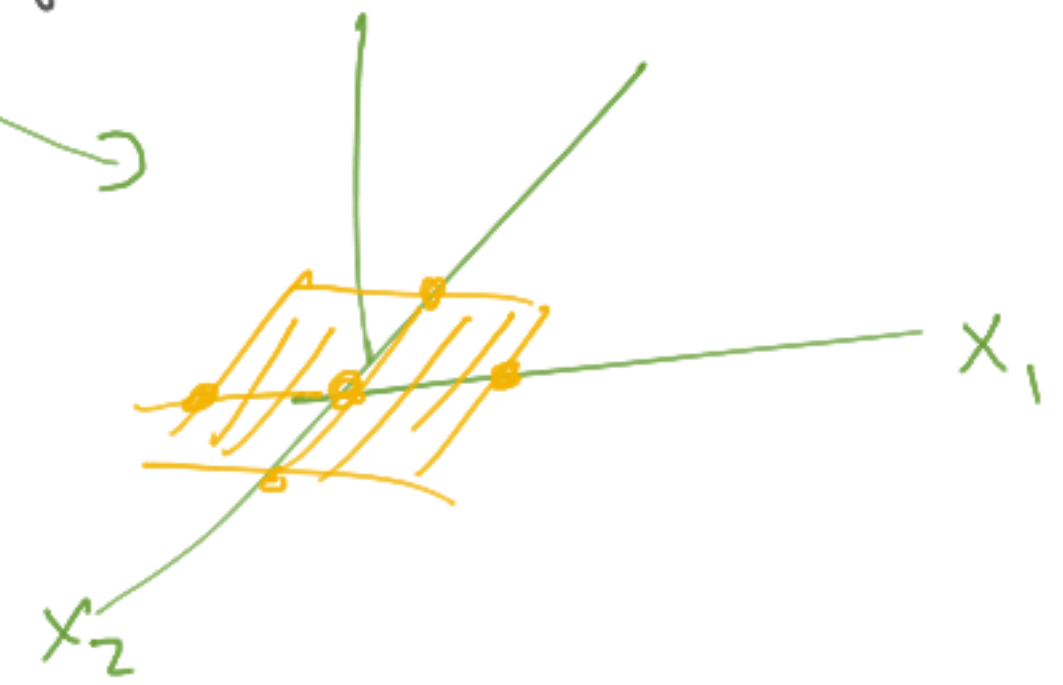


invariant to reparameterization.
 "Natural gradient"

same model



model space



f

g

$$x_{11} = 0$$

$$x_{12} = 1$$

$$y_1 = 1$$

$$x_{21} = 0$$

$$x_{22} = 0$$

$$y_2 = 0$$

→ one step
of ∇ descent
on k loss.