

Each generation

Starting pop

Ending pop

How many different individuals?

How much overlap?

How much functional diversity?

Can we use this as a measure of progression of learning?

Similar to early stopping

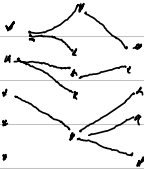
Day 11 presentation Aug 29 Tuesday Morning

Outline paper

outline presentation

30 minutes

writes down + lecture email



11 Friday of August
writes documents
Paper to him by Thursday
morning

Inter-~~island~~ migrations for reinforcement relative based interactions
optimization technique is not necessary tied to the
functionality of the architecture, but is tied
to its efficacy

Independent behaviors vs. goal or behaviors
(multi-objective behavior)
go to ROI avoid collision
vs.

some dynamic blend of
these two behaviors

What is actually meant by a "complex interaction"?

infl
like
obstacle migration } increase complexity
introduce collision penalties

Agent types

There are two types of agents acting in this environment:

each with a different purpose and dynamics in the environment. Harvesters gather resources, and excavators remove obstacles scattered throughout the environment. Both of these effects occur when the appropriate agent type collides with its resource. That is, when a harvester collides with a resource, it gathers the resource, ~~and the global reward is increased by the value of that resource~~ when an excavator collides with an obstacle, that obstacle is removed from the environment.

simplifying assumptions
ideal sensors
ideal actuators
homomorphic

However, if the wrong type of agent tries to collide with a resource or an obstacle then the actions of the agent is disordered, and the agent does not move.

Obstacles and rewards

Reward structure

In this work, we are primarily interested in how we can reinforce the correct behavior for goats. However, this may also entail negative reinforcement in the case of a detrimental action. Thus, we explored two reward structures as learning signals for the goats (these rewards are all global)

{ The first reward is purely positive reinforcement. Whenever a harvester collects a resource, the global reward is positive rewards for affecting the environment in a beneficial way

{ include penalties when goats collide w/ the way take at resource. Excavators do not want to destroy harvest resources and harvesters do not want to damage their machinery for colliding with an obstacle
deliberate

Experimental configurations

- algorithmic compensators
- scale restrict goats to obstacles
- scale zero of goats to resources
- scale reward penalties
- Transfer team learned in one ^{env} variation to another env variation (Montez)

Allow for resources to spawn different

- Monkeys also learn a position distribution for likelihood of resources being present at certain locations
- Bound locations of resources to particular areas
- need another type to scout out these areas

Ucla \rightarrow MFI

Island monkeys

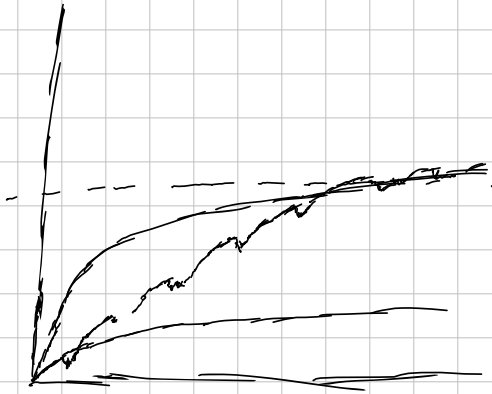
early visitors because...

introduction approach

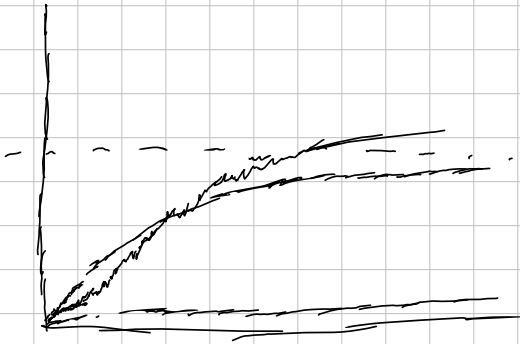
pro gen graphs

\rightarrow videos of behaviors
act in env

No collision penalty

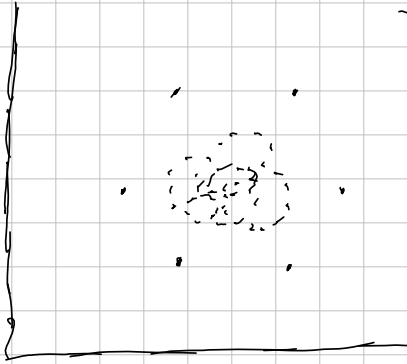


collision penalty



Keep in mind however
MFL is pretrained with
policies that are explicitly
chosen to be successful in
the environment. The Harvester
are trained to collect
resources and the excavators
are trained to remove obstacles,
each when set in the
presence of the other type
of environment features.

Show for both hermitages
and exheritages



- Neotomas at low common
- the goat is in that area
- with very limited movement
towards on
- less penalty, but less reward
- MPL more movement, but also
more interesting
- Instant balance - cost for
exaggeration before traversing
bad areas

look as MPL paper for
neotomas regions?

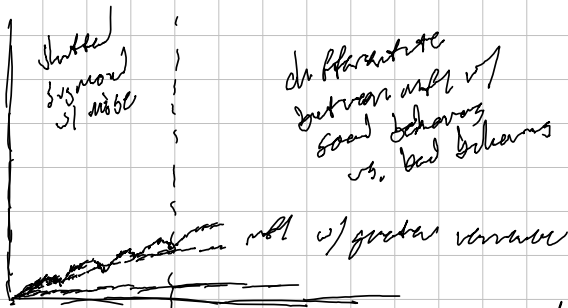
- Use simple order of
policy assignment
- might be getting hung up
on changes always being
in 1, 2, 3 positions
- run on generation instead of dist
- generation into own block
- check w/ Jupyter
- fine generation policy names

importance of
(non) based migration

- decreasing schedule heuristics
superior overall

Select names
room than
policy
directly

Pool of
policies



Happens less when there
is lower penalty for
collisions

How fast runs
→ didn't into a graph corner

✓ no penalty, will learn
an interesting policy due
it actually chooses to
go into negative regions

Clarify on issue of communication

Note 12

What is actually being shown?
(loss?)

Background slide 19

so what was everything prior?

Note 33

explain abstract

What do you mean by this?

Slide 47 What is the dog law here?

Verify matters credits

Challenges

describe as relevant

→ then state them

→ get into issues w/

counter examples

Struggle over the first
before moving on