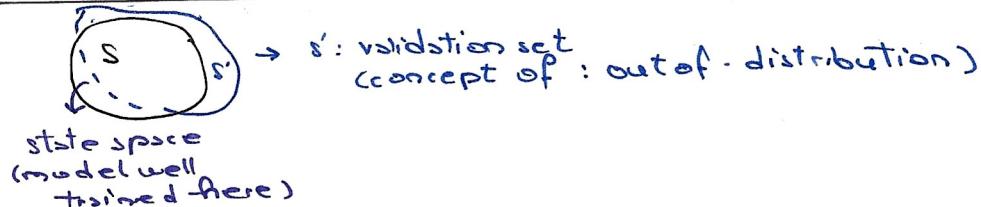


GENERALIZATION FOR ROBOTICS

What is Generalization?

- Can be working with diff. versions of the state space that hasn't seen before.
- Perform well on similar problems
- Out of distribution:



Why is Generalization Good?

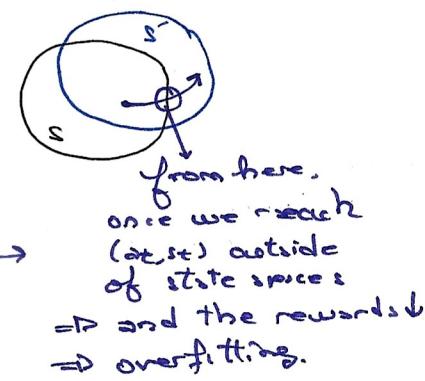
1. Do More with Less (information or structures)
 - The agent works even if it has not exp. that state. \rightarrow strong form of data efficiency
2. We generalize well because we make good use of data.
3. Good Agents should generalize to similar situations \rightarrow discover best in.d./causal mechanisms
 - Compression, less deep connections.
 - Occam's Razor
4. Avoid Overfitting. \Rightarrow Data Augmentation + adversarial examples solutions.

Policy $\Rightarrow \max IES^*(\pi)$ for state space
 \Rightarrow However if $IES^*(\pi)$ low: overfitting in RL
 \Rightarrow We want the policy to generalize on sequences of states?

[S: some computation simulation]
 [S': the real world]

$$\max \mathbb{E} \left[\sum_{t=0}^{\infty} r(s_t, a_t) \right]$$

$\sum_{s \sim p(s, a)} \mathbb{P}(s, a)$
 $s \leftarrow p_{\text{sim}}(s, a)$
 $s' \leftarrow p_{\text{real}}(s, a)$



Overfitting to the

1. Dynamics (distribution p(s,a)) \rightarrow \emptyset use hands after putting gloves
2. Rewards \rightarrow task misspecification
 \emptyset want to be a poker player for example.

5. Avoid Wasting Experience

Generalization in Robotics

Overfitting in...

- To a specific state distribution
- To the dynamics
- To Reward functions

Perceptual Generalization

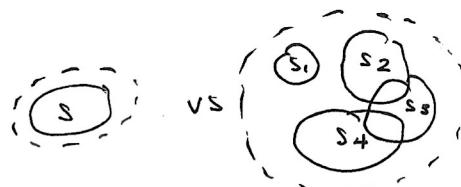
- Learning to overfit appearance + \Rightarrow with physics changes
- Method: Data Aug. + Adversarial Examples

Gen. in Dynamics

- $p(s_{t+1} | s_t, a_t)$: forward dynamics can change between tasks
- Ex: Motor down or breaks.
- Ex: changes in friction on \emptyset surfaces

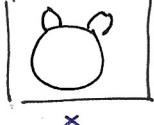
Multi-Task Learning

- We need the agent to train various tasks:



Transfer Learning

- Distribution $p(y|x)$, \exists lower dim. manifold (z) where $p(y|z) = p(y|x)$
- Want to learn transformation $p(z|x)$ \downarrow lower dimension

Ex:  \rightarrow $\begin{bmatrix} \text{has ears} \\ \text{nose} \\ \vdots \\ z \end{bmatrix}$

- {
- ① CNN Model from scratch
 \downarrow
Knowledge
 \downarrow
 - ② Pretrained CNN
- }

- If there's an overlap between tasks

- Training one improves perf. on others
- Learning a better $p(z|x)$ for future tasks.

Multi-Task Learning \Rightarrow Learn to be robust to tasks.

init θ to rand. network and $D \leftarrow \emptyset$

while True do

for $i \in \{0, \dots, n\}$ do

$T \sim p(T)$ \Rightarrow Select a task

for $t \in \{0, \dots, T\}$ do

Select + Apply Action $a_t = \pi(\cdot | s_t, \theta)$ in $p(s_{t+1} | s_t, a_t, T)$ \Rightarrow collect data for environment

Get experience $\{s_t, a_t, r_t, s_{t+1}\}$ and add to D . \Rightarrow collect exp. across env.

end for

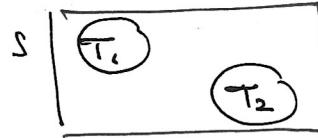
end for

Update policy θ using D , then $D \leftarrow \emptyset$

end while

Challenges

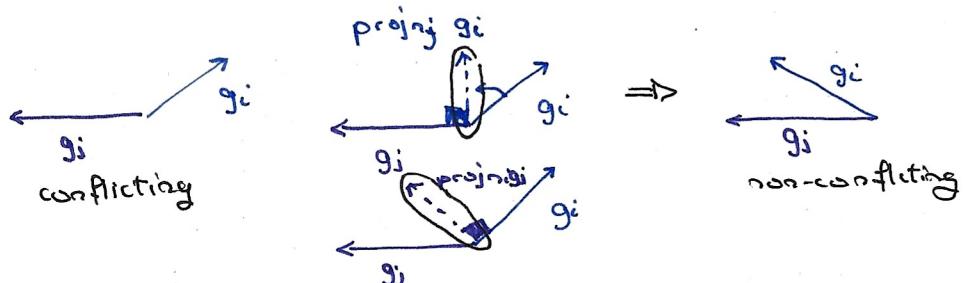
- Progress not easy
 - Goal: increase training speed + gen.
 - Multi-Task often makes things worse
- Cause: Interference



- The agent will focus on easiest Task (starts to overfit on some specific tasks)
 - Optimization difficulties
 - Asking policy to learn more
 - No free lunch theorem
- $p(z|x) \approx h(s)$

INTERFERENCE

- Gradients will contradict each other.
- Solution: PCGrad (Projection)



Focusing on Easiest Task

- Common issue in many methods (Options, HRL, Multi-Tasking)
 - Always check the video results.

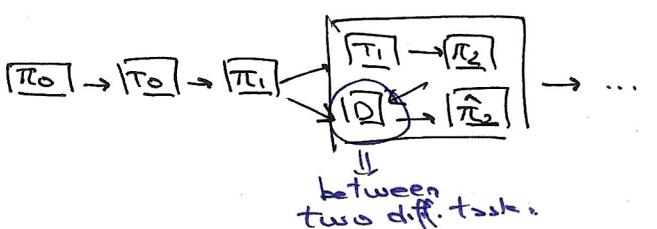
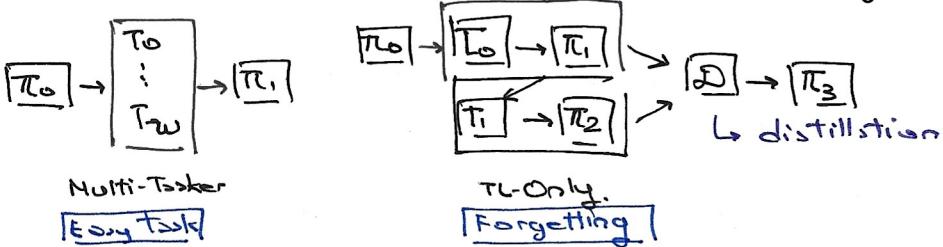
Methods:

1. Train over tasks separately \oplus combine using

behaviour cloning (Dagger)

↳ Combining new skills together.

↳ Distillation.
doesn't rely on reward functions.



PLAiD (better version and transferability than TL-Only)
→ Avoid forgetting.

→ We must «forgetting» and «quality»

↳ TL: difficult to predict effectiveness + transfer value function
+ overfitting specific tasks (\uparrow forgetting)

↳ PLAiD: \downarrow quality, but \downarrow forgetting

TL: \uparrow quality, but \uparrow forgetting

Meta Reinforcement Learning

but before that: Discussion about a paper

PROGRESSIVE RL WITH DISTILLATION FOR MULTI-SKILLED MOTION CONTROL

Combination of Distillation, Transfer Learning and MAML.

behavior cloning. (Dagger)

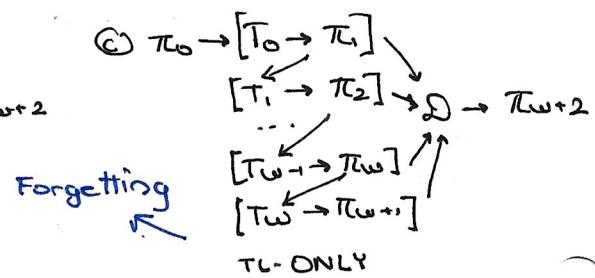
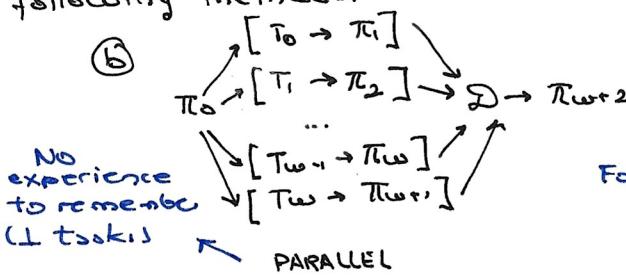
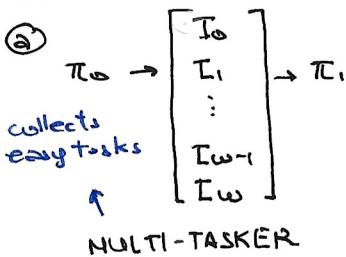
Transfer Learning

- Expert on source $T_i \Rightarrow$ post-training π_i
- Student on target task $T_{i+1} \Rightarrow \pi_{i+1}$, trying to match the expert's expectation.

Ex: 2D Biped Controller \Rightarrow Train on different types of terrains as tasks → incline, steps, slopes, gaps, etc.

Sequence Learning Schedules

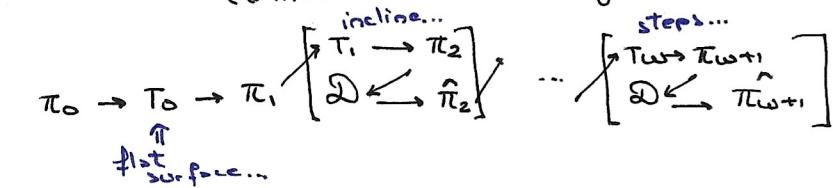
Issue with the 3 following methods:



*: Concept of «forgetting»

- Once we get to Task(w-1) for example, the T_0 will be totally forgotten
- Hence, if there was overlap between them, no more advantage on them anymore.

(d) PLAID: adding a new task each time + distillation, solves "forgetting"
(a little bit), better forward transfer abilities.



⇒ Generalization: incorporating new tasks while remembering the previous ones.

→ Measure effectiveness with forgetting and quality.

- TL: hard to predict effectiveness + overfitting specific tasks + forgetting
- Distillation needs some work for PLAID.

MetaRL: learn to learn

- so far, discussed how to generalize to unseen tasks
- What if we can collect a little data from the target task

$$p(z|x) \text{ vs. } p(z|x_{\text{AD}})$$

↑ ↑ ↴
RL Net-RL

- Collect a little bit more data on the target task and make a big jump in the reward function overall for that task.

[Ex:] Using exp from other Tasks.

- Board games (a new one) \Rightarrow you have played many before
 \Rightarrow similar dynamics \oplus figure out how to maximize pts.
- Meta-Learning: operationalize this process.

\curvearrowright

to predict if we can win the board game.

FRAMEWORK

- Given a function $f(x; \theta)$ with parameters θ and a loss function \mathcal{L} .
- Samples $D_i = (x_i, y_i)$ $\xrightarrow{\text{predict with loss based on the action/config. } x_i}$
 $\xrightarrow{\text{images/configurations of the board}}$
- \rightarrow Task-Loss: $\mathcal{L}(D_i; \theta_i) = \mathbb{E}_{(x_i, y_i) \sim D_i} [\mathcal{L}(f(x_i; \theta_i), y_i)]$
- In supervised setting: use paired data of input x_i and output y_i
- Goal: Enable fast adaptation on a diff. set of meta-test tasks not seen during training.
 \rightarrow How to win frequently in this new board game with a few data from this board?

\Rightarrow Solution: MAML Algorithm

\curvearrowright

MAML: Model-Agnostic-Meta-Learning. \Rightarrow LIBRARY: **HIGHER**

- Fixed distribution of tasks $p(T)$ \Rightarrow all the board games in the shelf.
- \rightarrow M tasks drawn from $\{T_i\}_{i=0}^M$ distribution
- \rightarrow Tasks $T_j \sim p(T)$ that provided data $D_j = \{x_j, y_j\} \Rightarrow D_j$: play games for board game $\langle j \rangle$
(assumes with some ppl)
- \rightarrow Train θ for a model on $\{T_i\}_{i=0}^M$
- \rightarrow Objective $f(x; \theta)$ is low after few gradient updates on data D_j

$$\min_{\theta} \sum_{T_i \sim p(T)} \mathcal{L}((\theta - \alpha \nabla_{\theta} \mathcal{L}(\theta, D_i^{\text{train}}), D_i^{\text{test}}))$$

$$= \min_{\theta} \sum_{T_i \sim p(T)} \mathcal{L}(\phi, D_{T_i}^{\text{val}})$$

OTHER METHODS.

Train RNN to ingest description and predict a MDP

- Sequence of data that helps you to predict a policy that will max. rewards inside this env.
- Overfitting

Train optimizer or update function **HAML** $\rightarrow p(z|x, D, \theta)$

- Good extrapolation + consistent
- Hard to determine # samples or gradient steps

Train a representation $p(z|x)$ **PEARL**

- Inference problem
- Overfitting (doesn't generalize well)

Transformers: works well for diff. policy changes (humans for example)

GENERALIZATION WITHOUT TASKS.

1. Learn to infer the boundaries.
2. Learn without explicit task boundaries (detect local state distributions)

LIMITATIONS.

- Meta-learning can memorize.
 - Pickup on obvious indications of what the task is.
 - ⚡ for new data
 - Need to be a good Bayesian + avoid assuming model is good.

Need good average policy

- Training process: very data-sufficient. (lots of data)