# Adaptive spatial discretization using reinforcement learning

Jemil Avers Butt[1,2], Andreas Wieser[1]

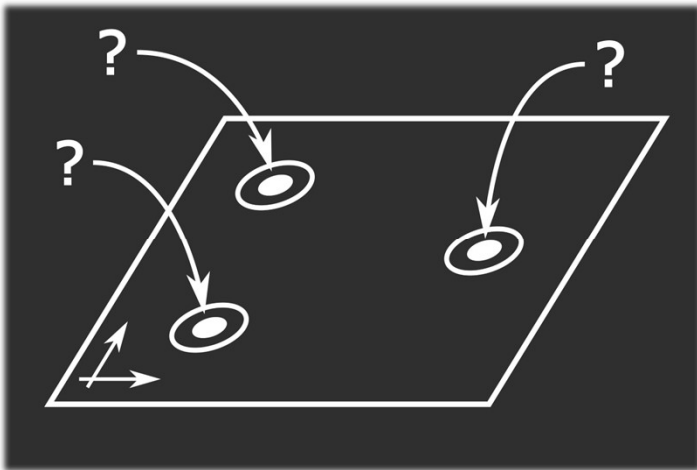[1]ETH Zurich, Institute of Geodesy and Photogrammetry, Zurich;  [2]Atlas optimization GmbH, Zurich
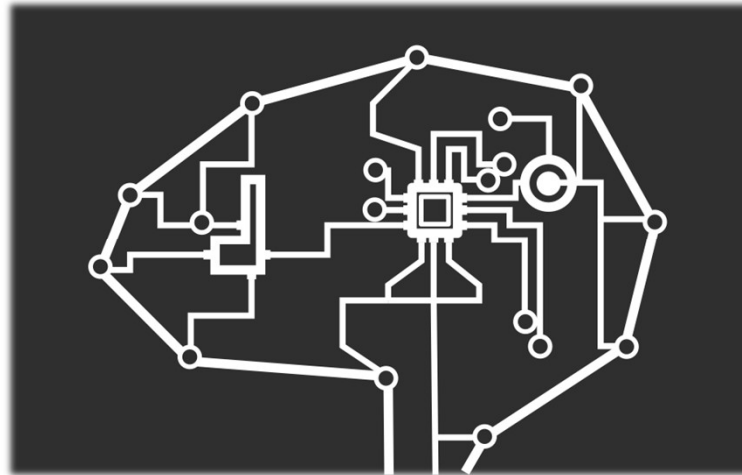
**ETH**zürich

ATLAS
OPTIMIZATION

# Content

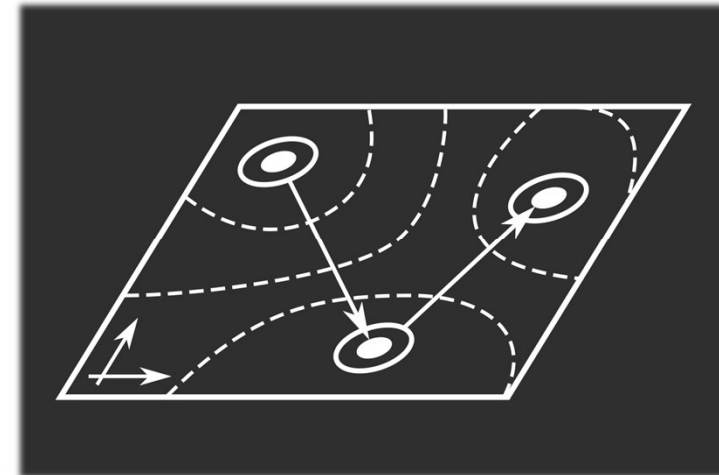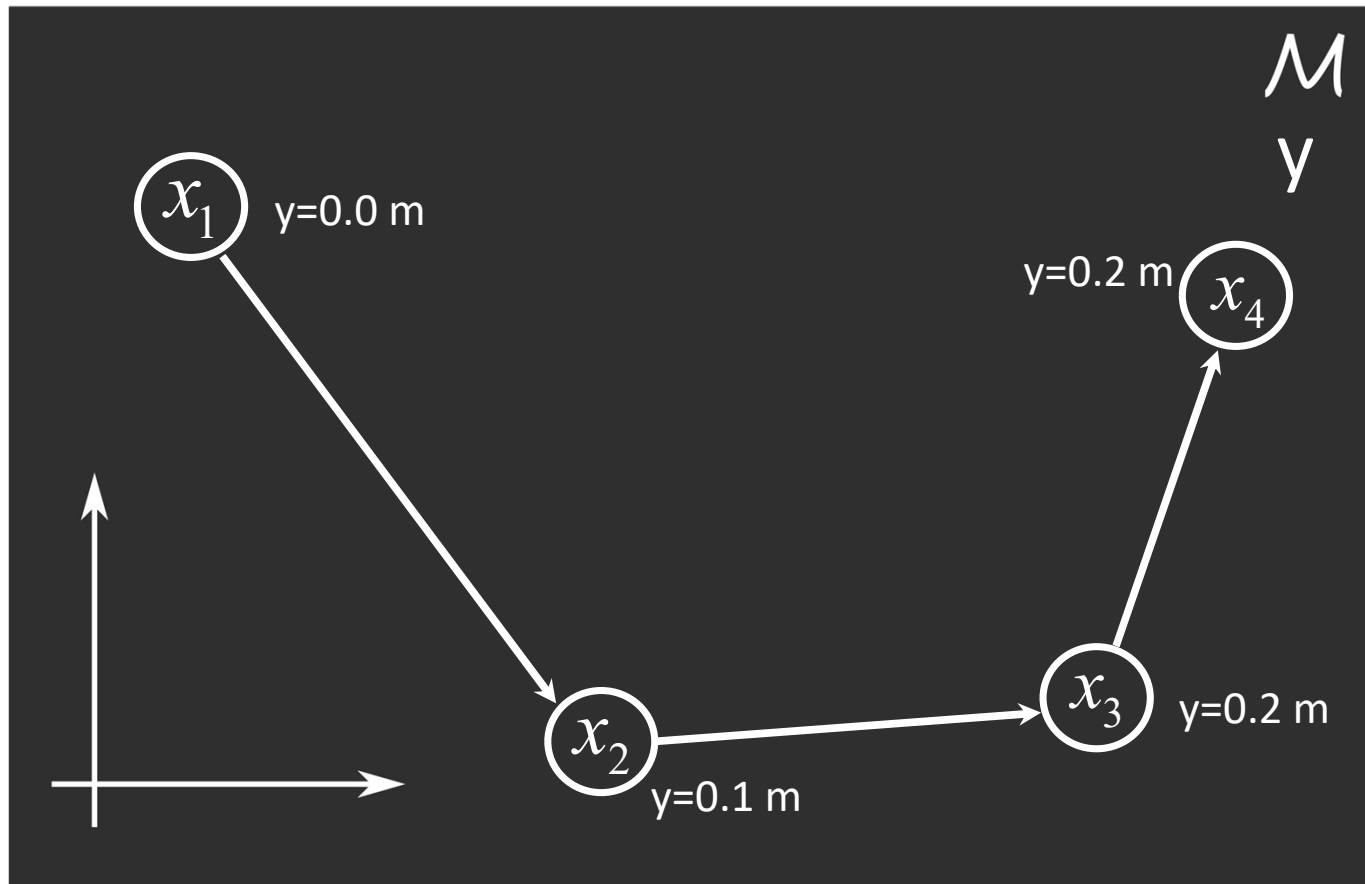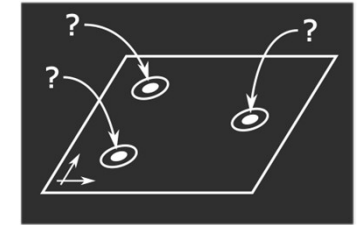| 1. Problem | 2. Method | 3. Results |
|:---:|:---:|:---:|
| Spatial discretization | Reinforcement learning | Better sampling |

# Problem - Specification



**Spatial discretization**

- Object $M$
- Interested in $y$
- Sequence of $x_1, x_2, \ldots$

**Questions**

- Best sequence?
- System knowledge?
- Adapt to observations?

**Answer**

- We need a policy $\pi$ !

# Problem – Hierarchies of Solutions

$\pi(\cdot)$ is called policy

$\pi : \text{state} \mapsto \text{action}$

What we know      Next measurement

## We want

- Adaptivity

  Next sample locations influenced

  by pevious findings

  $\rightarrow \pi = f(x_1, ..., x_N, y_1, ..., y_N)$

- Prior knowledge

  Sampling scheme makes use of

  dependences between $x, y$

  $\rightarrow \pi = f(x_1, ..., x_N; p_{xy})$

# Problem – State of the art



$$\pi = \pi(x_1,...,x_N)$$
"blind search"

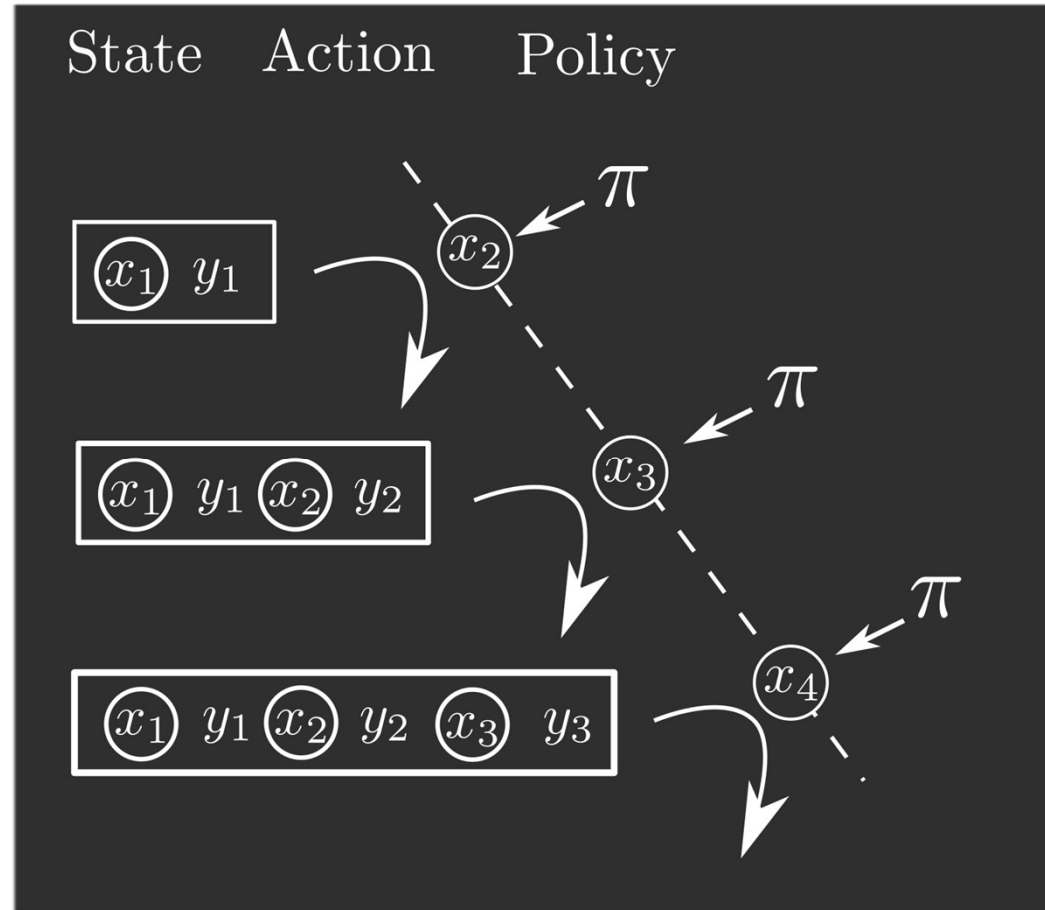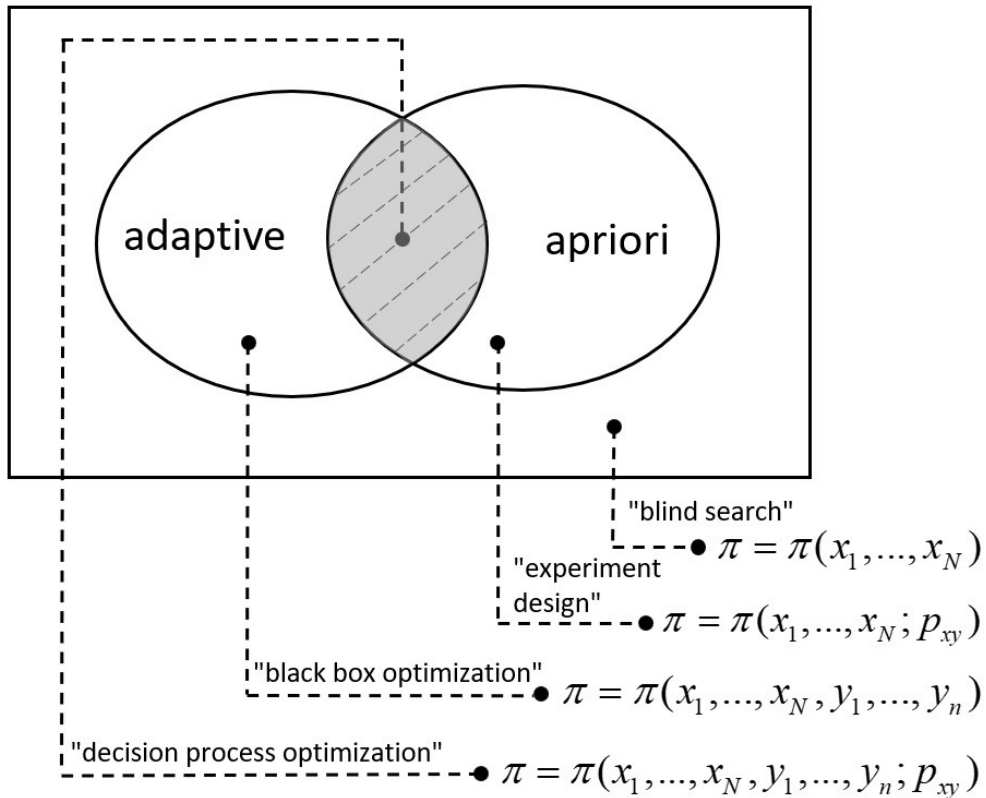$$\pi = \pi(x_1,...,x_N; p_{xy})$$
"experiment design"

$$\pi = \pi(x_1,...,x_N, y_1,...,y_n)$$
"black box optimization"

$$\pi = \pi(x_1,...,x_N, y_1,...,y_n; p_{xy})$$
"decision process optimization"

"Blind search" :
>   Random numbers, pseudorandom, grids, ...
>
>   [Kuipers & Niederreiter]

"Experiment design" :
>   Minimization of $\mathrm{tr}(\Sigma),\ \det(\Sigma),\ \mathrm{Entropy},...$
>   - > Geodetic networks, geostatistical sampling
>
>   [Grafarend & Sanso], [Angulo et al]

"Black box optimization" :
>   Sequences of samples based assumptions.
>   -> Problems with unknown structure
>
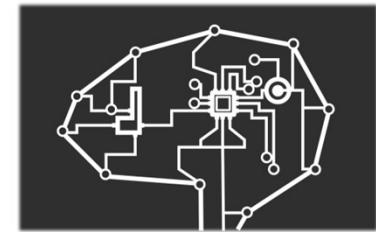>   [Fu], [Alarie et al]

"Decision process optimization" :
>   Adaptive policies to maximize reward
>   -> Optimal control in economy and games
>
>   [Sutton & Barto],
>   [Feinberg & Shwartz]

# Method – RL in ML

Clustering

Supervision

Difficulty

Regression

Dimensionality
reduction

ML Problems

Classification

Reinforcement learning

RL:
- Optimal control

- No right/wrong
- Learn from
  experience
- Agent interacts
  with environment

- Difficult task
- Computationally
  intense

# Method – RL Policy gradients

Policy function is ANN
- Determines decisions
- Make decisions optimal

1. Expected return depends on params

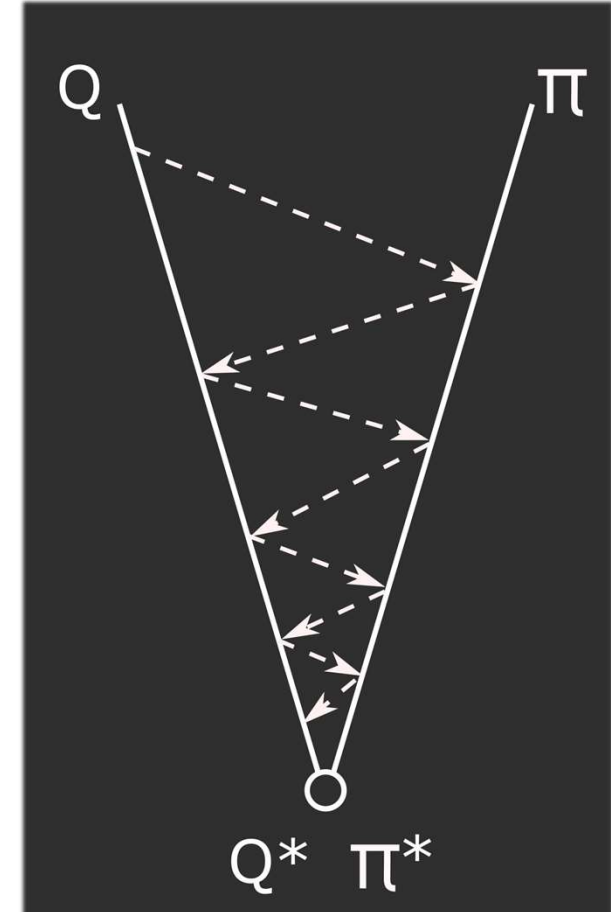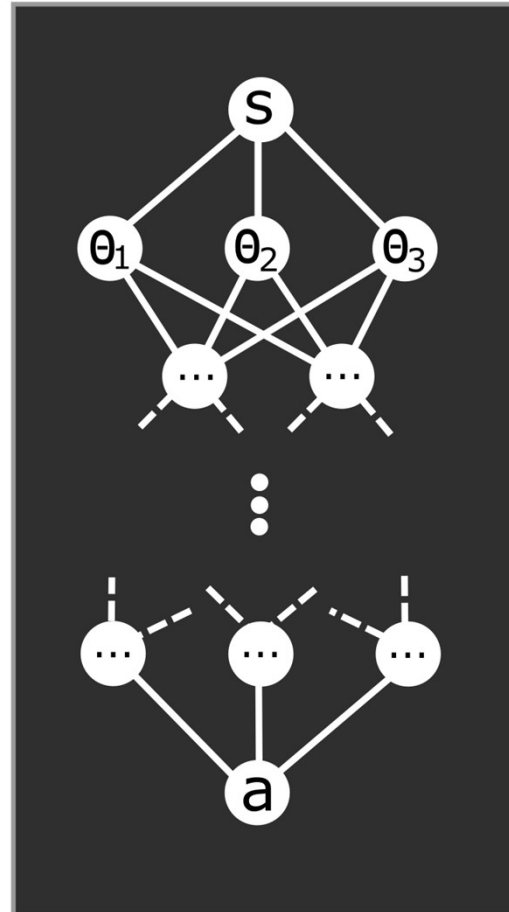$$\eta(\pi) = E_{\pi p}\left[\sum_{k=1}^{n-1}\gamma^k r(s_k, a_k)\right]$$

$$\eta(\pi) = \int_{(S\times A)^{n+1}}\sum_{k=0}^{n-1}\gamma^k r(s_k, a_k)p_\theta(\xi)d\xi$$

2. Gradient depends on observations

$$\nabla_\theta\eta = \int_{(S\times A)^{n+1}}\sum_{k=0}^{n-1}\gamma^k r(s_k, a_k)\nabla_\theta[\log p_\theta(\xi)]p_\theta(\xi)d\xi$$

$$\hat{\nabla}_\theta\eta = \frac{1}{m}\sum_{j=1}^{m}\left(\sum_{k=0}^{n-1}\gamma^k r(s_k^j, a_k^j)\right)\nabla_\theta[\log p_\theta(\xi_j)]$$

$$= \frac{1}{m}\sum_{j=1}^{m}\left(\sum_{k=0}^{n-1}\gamma^k r(s_k^j, a_k^j)\right)\sum_{i=0}^{n-1}\nabla_\theta\log\pi_\theta(s_i^j, a_i^j).$$

# Method – Showcase

Training: 3 min

Training: 3 h

Training: 3 d



Trained with a3c by David Griffis
https://github.com/dgriff777/a3c_continuous

# Method – Spatial discretization as RL



Train: $\dot{\uparrow} + \uparrow$

Test: $\uparrow$

- RL algorithm
- Environment

Action $\pi(s_t) = x_{next} = a_t \in A$

$x_{next}$

$\pi$

RL

$\pi_{new}$

State $s_t$

State $s_t \in S$

$t \rightarrow t+1$

$x_{next}$

New state $s_{t+1}$

Reward $r_t = r_t(s_t, s_{t+1}, a_t)$

Spyder (Python 3.8)

File  Edit  Search  Source  Run  Debug  Consoles  Projects  Tools  View  Help

/home/jemil/Desktop/Programming/Python/Atlas_Optimization/Optimal_Monitoring/Compilation_Paper

...rogramming/Python/Atlas_Optimization/Optimal_Monitoring/Compilation_Paper/benchmark_random_def_2D.py

random_def_1D.py    Playing_games_Lunar_lander_discrete.py    Playing_games_bipedal_walker.py    benchmark_random_def_2D.py

```python
"""
The goal of this script is to train a TD3 RL algorithm on the random deformation
task and compare the cumulative rewards to the ones gathered by alternative
discretization strategies.
For this, do the following
    1. Definitions and imports
    2. Train with stable baselines
    3. Apply alternative methods
    4. Summarize and plot results
"""


"""
    1. Definitions and imports
"""


# i) Import basics and custom environment

import numpy as np
import time
from scipy.optimize import basinhopping
import class_random_def_2D_env as def_2D


# ii) Import stable baselines

from stable_baselines3 import TD3
from stable_baselines3.common.env_checker import check_env


# iii) Initialize and check

np.random.seed(0)
def_2D_env=def_2D.Env()
def_2D_env.reset()
# check_env(def_2D_env)



"""
    2. Train with stable baselines
"""


# i) Train a TD3 Model

# start_time=time.time()
# model = TD3("MlpPolicy", def_2D_env,verbose=1, seed=0)
# model.learn(total_timesteps=100000)
# end_time=time.time()

# model.save('./Saved_models/trained_benchmark_random_def_2D')
model=TD3.load('./Saved_models/trained_benchmark_random_def_2D')


"""
    3. Apply alternative methods
"""


# Note: All actions are in [-1,1]x[-1,1] and get mapped to [0,1]x[0,1] by
# the environment translating input actions from the symmetric box space
# [-1,1]x[-1,1] to indices

# i) Grid based sampling

def grid_based_sampling(environment):
    grid_x1=np.kron(np.array([-1/3, 1/3, 1]),np.array([1, 1, 1]))
    grid_x2=np.kron(np.array([1, 1, 1]), np.array([-1/3, 1/3, 1]))
    grid=np.vstack((grid_x1, grid_x2))
    action=grid[:,environment.epoch]
    return action


# ii) Pseudo random sampling

def pseudo_random_sampling(environment):
    Halton_sequence=np.array([[1/2, 1/4, 3/4, 1/8, 5/8, 3/8, 7/8, 1/16, 9/17,
```

Training reward



Console 1/A

```
[ 0.07423399]
 [ 0.36265202]
 [ 0.28539754]
 [ 0.35397827]
 [-0.02127369]]
Reward is  -0.0012474826886026413
Measured locations are [[0.20408163 0.87179487]
 [0.55102041 0.28205128]
 [0.3877551  0.64102564]
 [0.51020408 0.64102564]
 [0.28571429 0.30769231]
 [0.65306122 0.51282051]
 [0.71428571 0.64102564]
 [0.42857143 0.92307692]
 [0.         0.        ]]
Measurements are [[ 0.27411945]
 [ 0.27618975]
 [ 0.35367949]
 [ 0.40609153]
 [ 0.07423399]
 [ 0.36265202]
 [ 0.28539754]
 [ 0.35397827]
 [-0.02127369]]
Reward means of different methods
[-0.04441617 -0.05062521 -0.08014935 -0.03963356 -0.04539691 -0.03088655]
Reward standard deviations of different methods
[0.04121231 0.03924319 0.06256735 0.03201321 0.03924964 0.0434375 ]


Figures now render in the Plots pane by default. To make them also appear inline in the Console, uncheck "Mute Inline Plotting" under the Plots pane options menu.


In [6]: n_episodes=1
   ...: for k in range(n_episodes):
   ...:     done=False
```
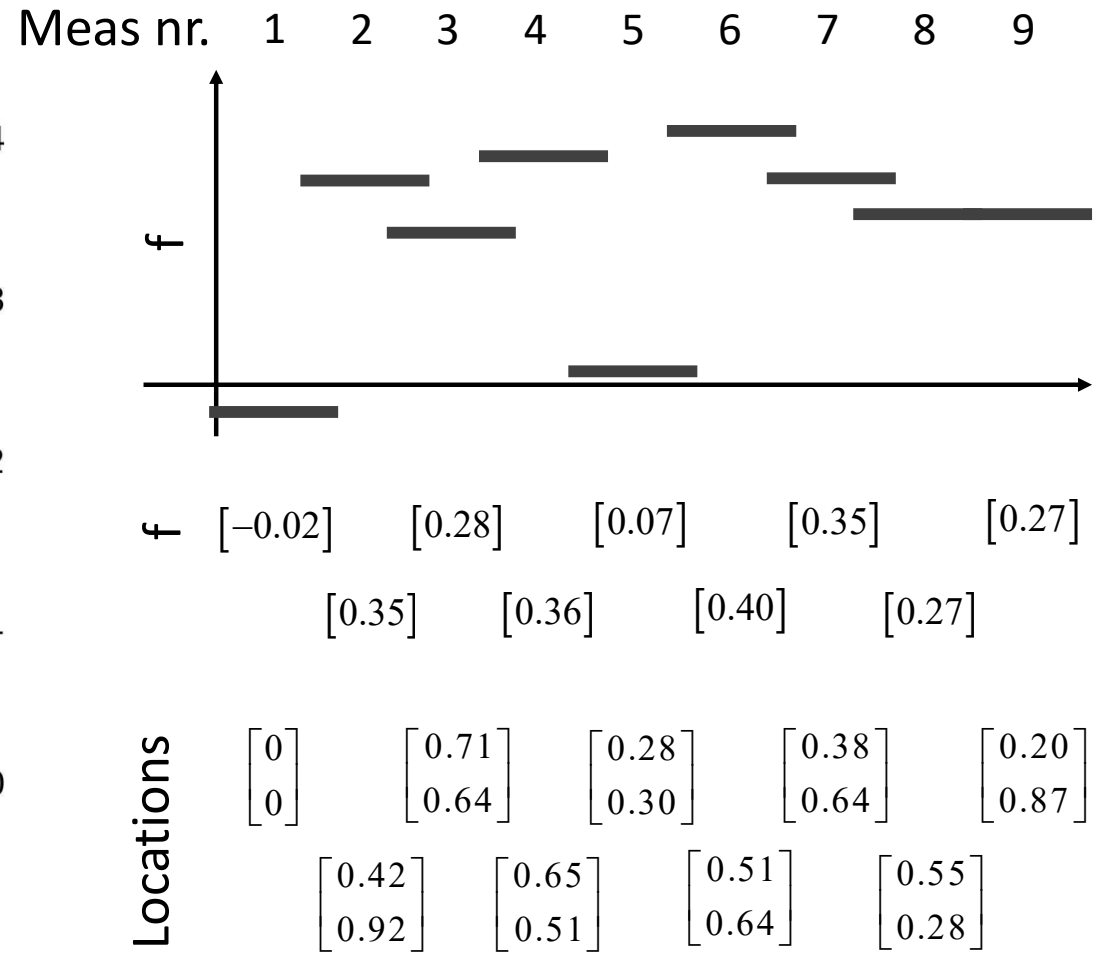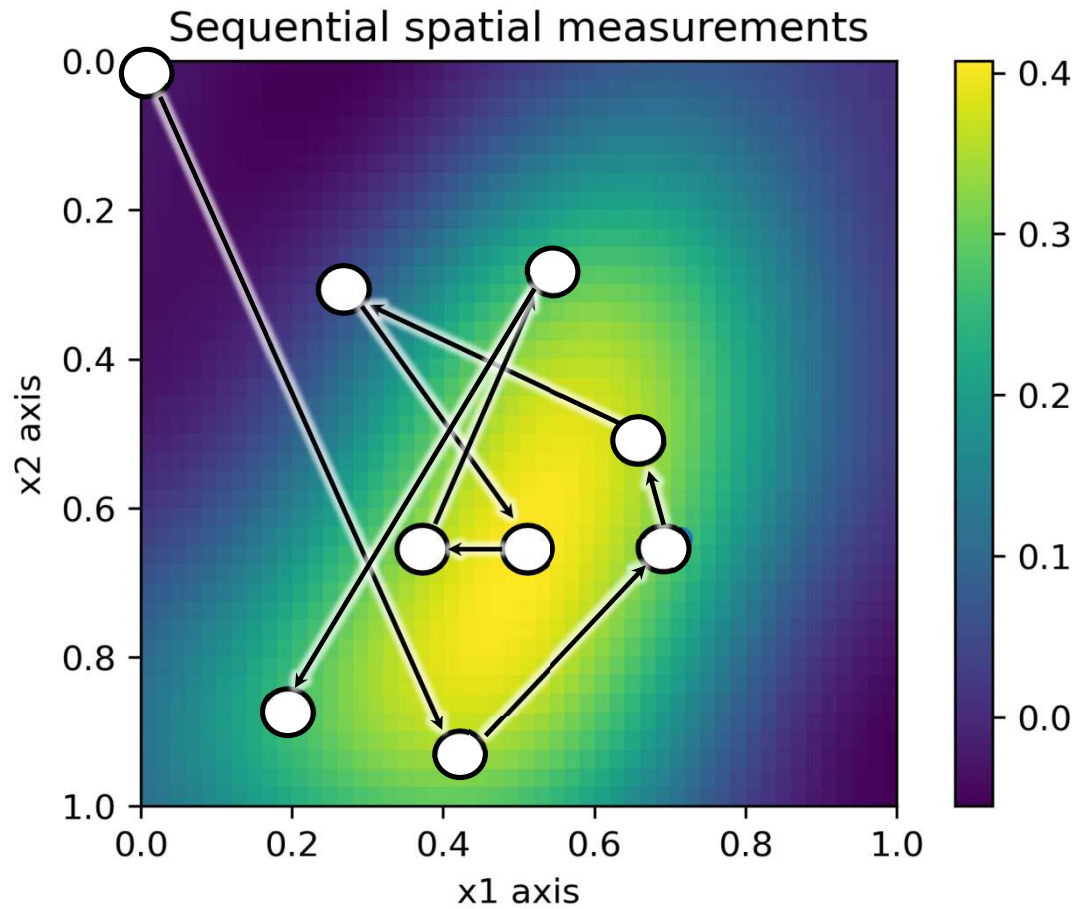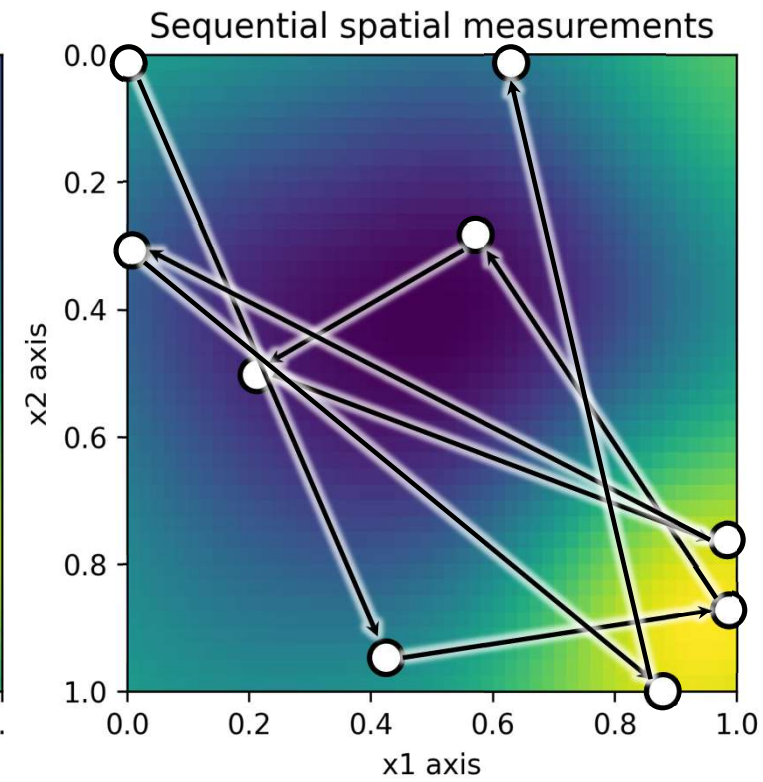
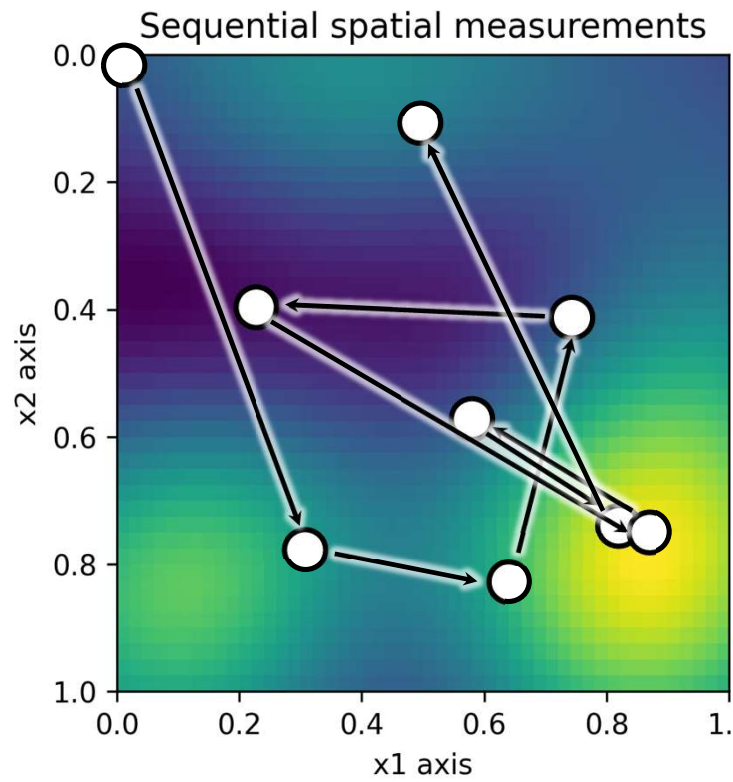Variable explorer  Help  Profiler  Code Analysis  Plots  Files

IPython console  History

LSP Python: ready      conda: base (Python 3.8.3)      Line 238, Col 17      ASCII      LF      RW      Mem 49%

# Results – Illustration



Sequential spatial measurements

# Results – Illustration



Sequential spatial measurements

# Results – Comparison



| Method/Task | Random | Quadrature | Exp. design | RL |
|---|---|---|---|---|
| beam bending | -0.05 ±0.1 | -0.065 ±0.07 | -0.0025 ±0.005 | **-0.0006 ±0.001** |
| 1D def. | -0.07 ±0.08 | -0.047 ±0.04 | -0.035 ±0.04 | **-0.019 ±0.03** |
| Def. tracking | -0.28 ±0.19 | N. A | -0.24 ±0.14 | **-0.067 ±0.038** |
| 2D def. | -0.08 ±0.063 | -0.04 ±0.032 | -0.045 ±0.039 | **-0.031 ±0.043** |

| Human* | Human 1 | Human 2 | Human 3 |
|---|---|---|---|
| 2D def. | -0.078 | -0.015 | -0.026 |

# Results – Conclusion

Problem:  Have some object or process. Want to know something.
Decide where to measure.

Solution:  Make model of process. Use it for simulations.
RL gets optimal* sequences of decisions.
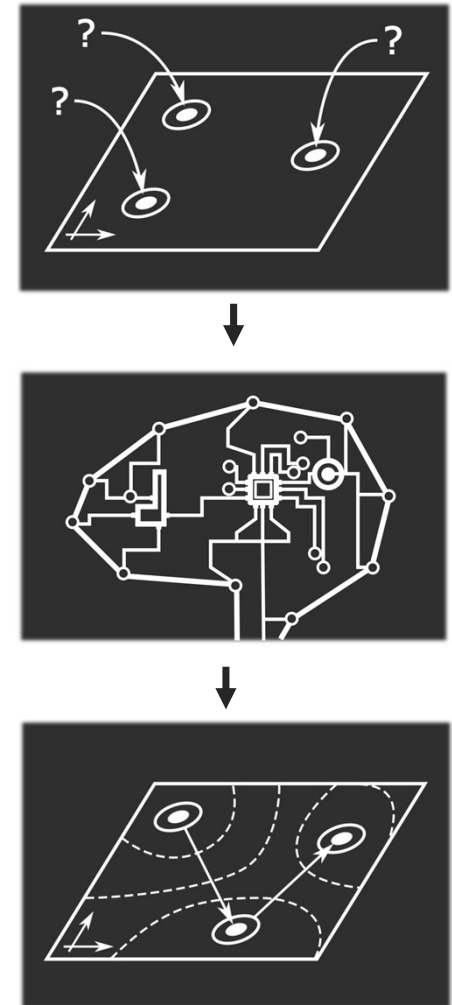
What's good!

- Really flexible
- Optimal decisions under
uncertainty
- Don't need explicit
probabilistic models
- Modelling work done
by computer
- Higher score than alternatives

What's bad!

- No guarantees
- No measure of uncertainty
- Need model for simulation
- Training sloooow …

Applications!

- Adaptive monitoring
- Episodic repositioning

# Sources

Ai, B., Jia, M., Xu, H., Xu, J., Wen, Z., Li, B. and D. Zhang (2021). Coverage path planning for maritime search and rescue using reinforcement learning, Ocean Engineering, Volume 241.

Alarie, S., Audet, C., Gheribi, A.E., Kokkolaras, M., and S. Le Digabel. (2021). Two decades of blackbox optimization applications, EURO Journal on Computational Optimization, 9.

Angulo, J.M., Ruiz-Medina, M.D., Alonso, F.J. and M.C. Bueso (2005). Generalized approaches to spatial sampling design. Environmetrics, Vol 16, No. 5.

Feinberg, E. A. and A. Shwartz (2012). Handbook of Markov Decision Processes. Methods and Applications. Berlin Heidelberg (Springer Science & Business Media).

Fu, M. C. Handbook of Simulation Optimization. Berlin, Heidelberg: Springer, 2014.

Grafarend, Erik W. and F. Sansò (2012). Optimization and Design of Geodetic Networks. Berlin Heidelberg (Springer Science & Business Media).

Sutton, R.S. and A.G. Barto (1992). Reinforcement Learning. Berlin Heidelberg (Springer Science & Business Media).

Wiering, M. and M. Otterlo (2012). Reinforcement Learning. State-of-the-Art. Berlin Heidelberg (Springer Science & Business Media).

Zuluaga, J. G. C., Leidig, J. P., Trefftz, C., and G. Wolffe, (2021). Deep Reinforcement Learning for Autonomous Search and Rescue, *NAECON 2018 - IEEE National Aerospace and Electronics Conference*, pp. 521-524.

# Thank you for your time!

## Now:

### Q & A



Sequential spatial measurements