

TExplore

Todd Hester and Peter Stone



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Structure

1. Introduction
2. Related Work
3. Background
4. Algorithm
5. Experiments

Goal

- Algorithm for real time targeted Exploration in large domains

General idea

- intelligent Exploration
- Sample efficient algorithm, which tries to explore a minimal number of states necessary to learn a good policy

Related Work

- There are a number of algorithms that address the exploration problem
- A few examples are R-MAX, SPITI, Bayesian RL methods or the Gaussian Process
- All have drawbacks, why they can't explore large domains in real time

Background – Markov Decision Process

- set of States S
- set of Actions A
- reward function $R(s,a)$
- transition function $P(s'|s,a)$
- discrete state is represented by a vector of n discrete variables

$$s = \langle x_1, x_2, \dots, x_n \rangle$$

Background – Bellman equation

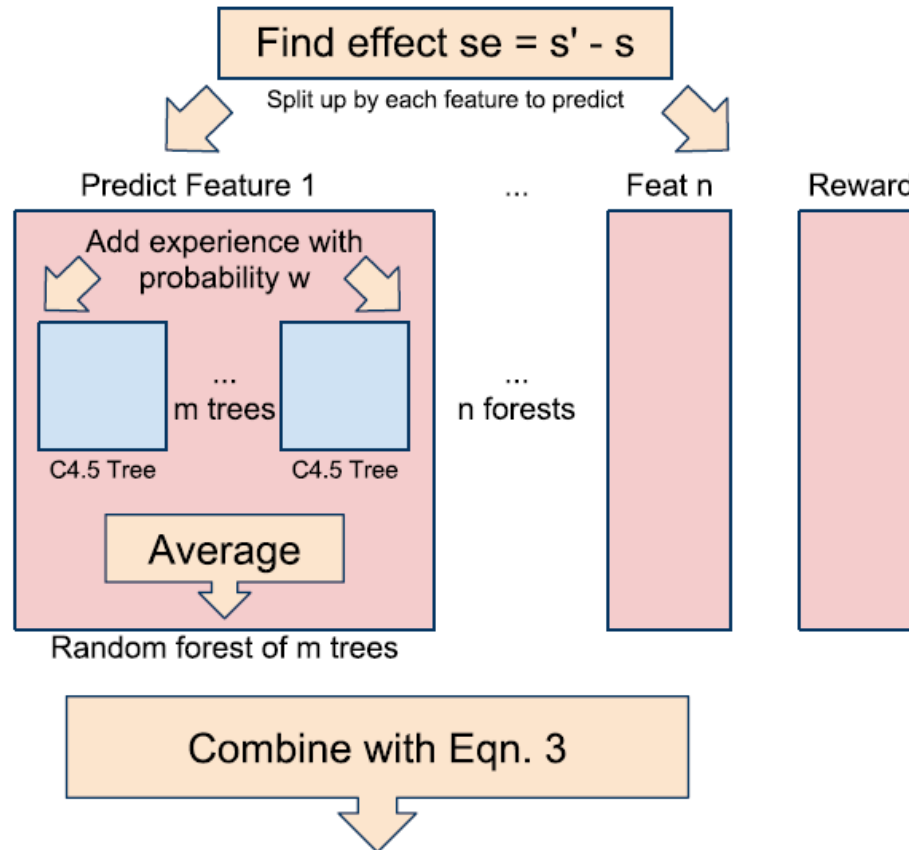
- $Q^*(s, a)$ of a state-action pair is an estimate of the future reward that can be obtained from (s, a)
- $Q^*(s, a)$ is determined by solving the Bellman equation :
$$Q^*(s, a) = R(s, a) + \gamma P(s' | s, a) \max_{a'} Q^*(s', a')$$
- Optimal policy $\pi(s) = \arg \max_a Q^*(s, a)$

- Easier to generalize the transition effect (relative change) than the exact outcome across states $\rightarrow s_e = s' - s$
- Learn separate prediction of each of the n state features and put them together to a complete model
- The probability of the transition effect is the product of the probabilities of its n state features : $P(s_e | s, a) = \prod_{i=0}^n P(x_e^i | s, a)$

Algorithm – Modul Learning

- assumes that each state feature can be predicted independently
- Each model is build by using a random forest
 - each decision tree is trained on a subset of the agent experiences
- Final prediction is the average of all tree outcomes
- Each tree is build recursively using the C4.5 algorithm

Algorithm – Modul Learning



- Agent is given intrinsic rewards based on the variance of the models predictions

$$R(s, a) = R_o(s, a) + b\sigma^2(s, a)$$

$$\sigma^2(s, a) = \frac{1}{n+1} \left[\sigma^2 R(s, a) + \sum_{i=1}^n \sigma^2 P(s_i^t | s, a) \right]$$

→ $R_o(s, a)$ is the original prediction of the model

→ $\sigma^2(s, a)$ is the variance in the model's predictions

Algorithm - Exploration

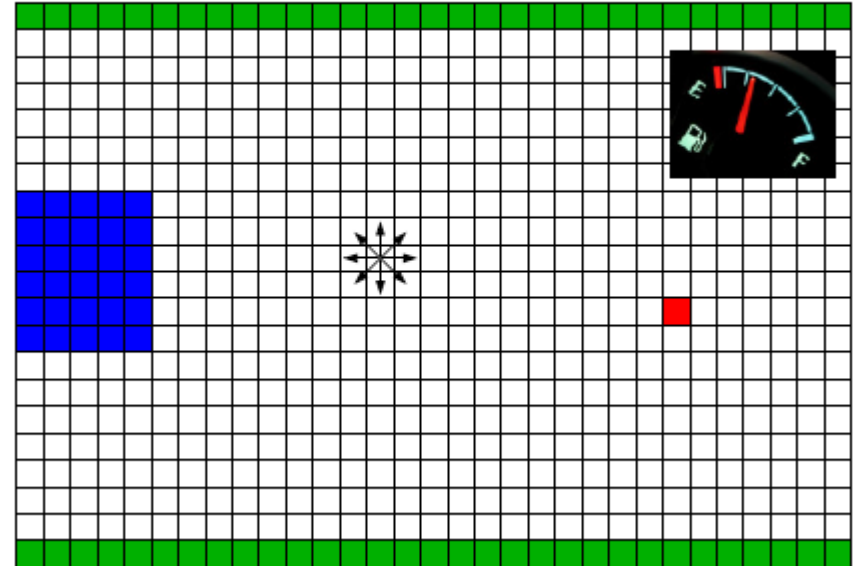
$$R(s, a) = R_0(s, a) + b\sigma^2(s, a)$$

- b is a coefficient either adding or subtracting intrinsic rewards from the model based on the variance
- $b < 0$, the agent will avoid uncertain states
- $b > 0$, the agent being driven to explore uncertain states
- $b = 0$, the agent will act greedily with respect to its model

- Algorithm needs to replan what the best action is, by using the UTC algorithm
- Provide basic knowledge of the structure of the domain
- Seeding the agent with a few sample experiences from the domain

Experiments - Environment

- 21 x 31 cells
- 61 possible energy levels
- 8 possible actions
 - 317.688 state-actions
- Agent fuel level between 0 and 60
- State vector s is made up to three features: Row, Col and Fuel
- Fuel station reward $R(x) = base - (x \bmod 5)a$
 - x = column, $base$ = reward for that row, a = variation of the costs across the columns



Experiments - Environment

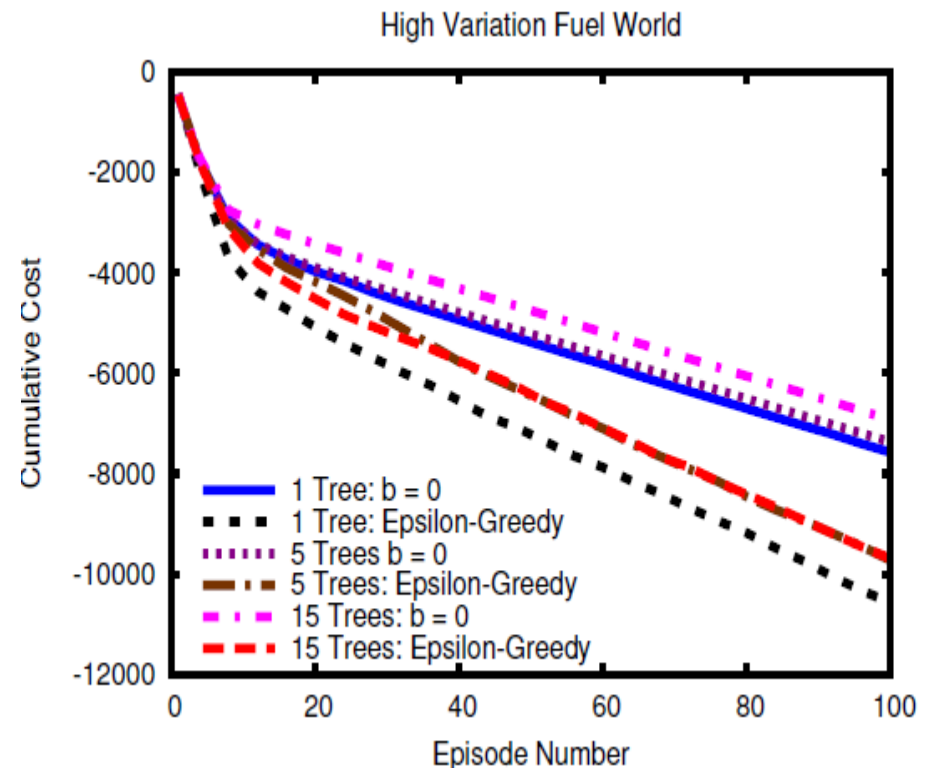
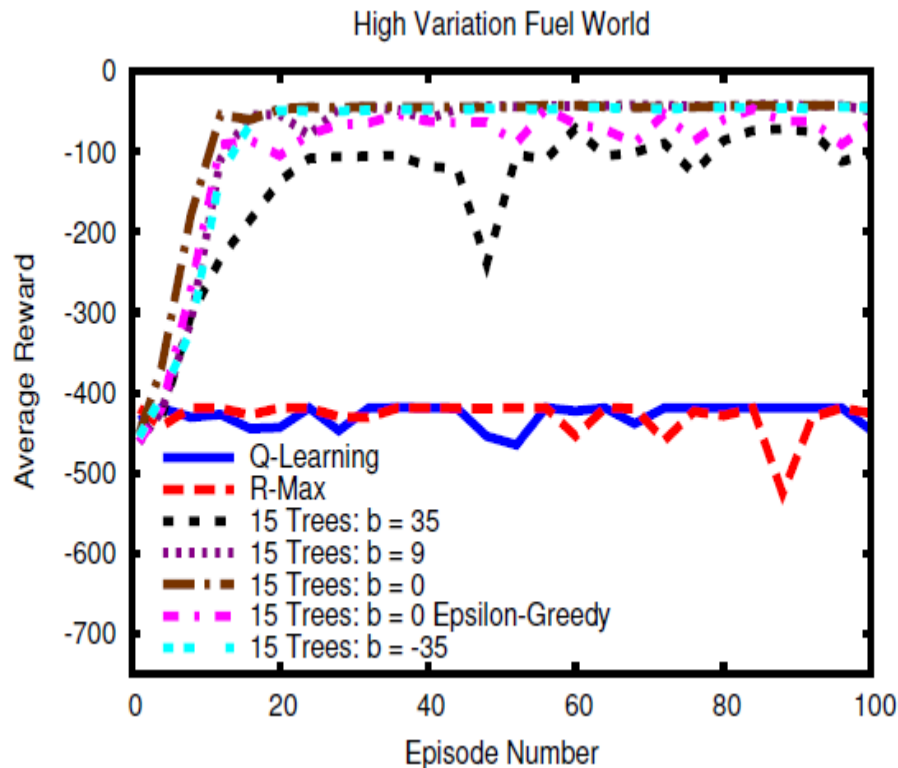
Low variation fuel world :

➤ $a = 1.0$, base = -18.0/-21.0 bottom/top row → reward = -18.0 to -25.0

High variation fuel world :

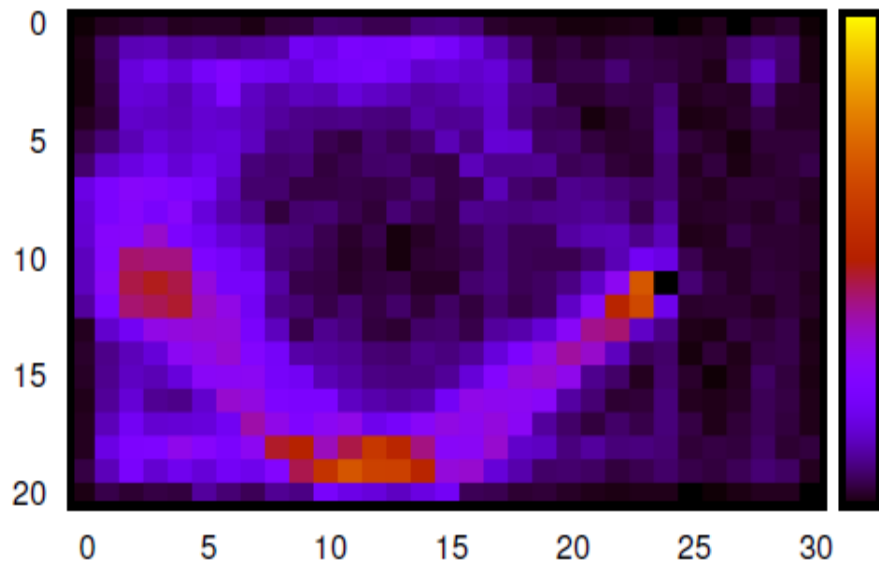
➤ $a = 5.0$, base = -10.0/-23.0 bottom/top row → reward = -10.0 to -33.0

Experiments - Results

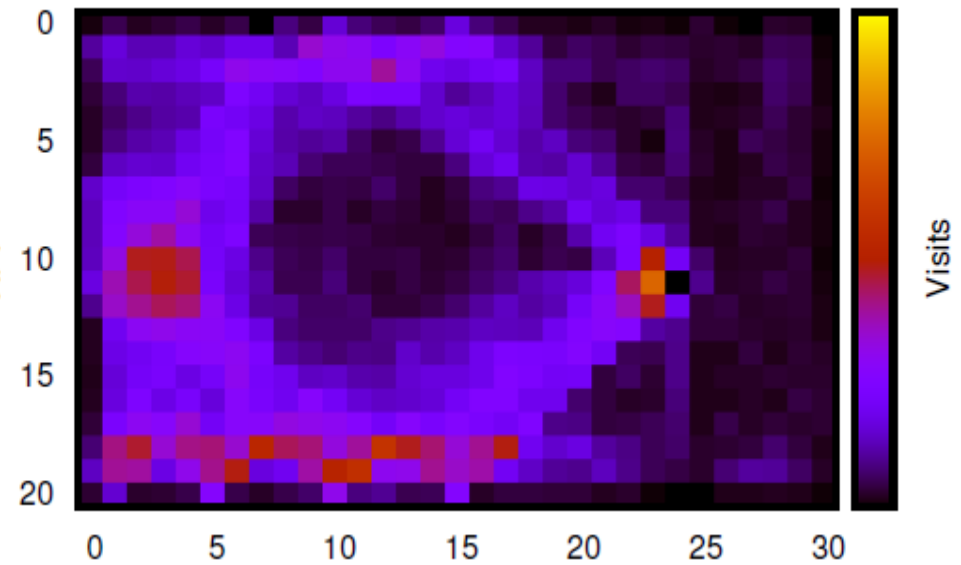


Experiments - Results

Low Variation Fuel World - First 20 Ep. - 15 Trees: $b = 0$

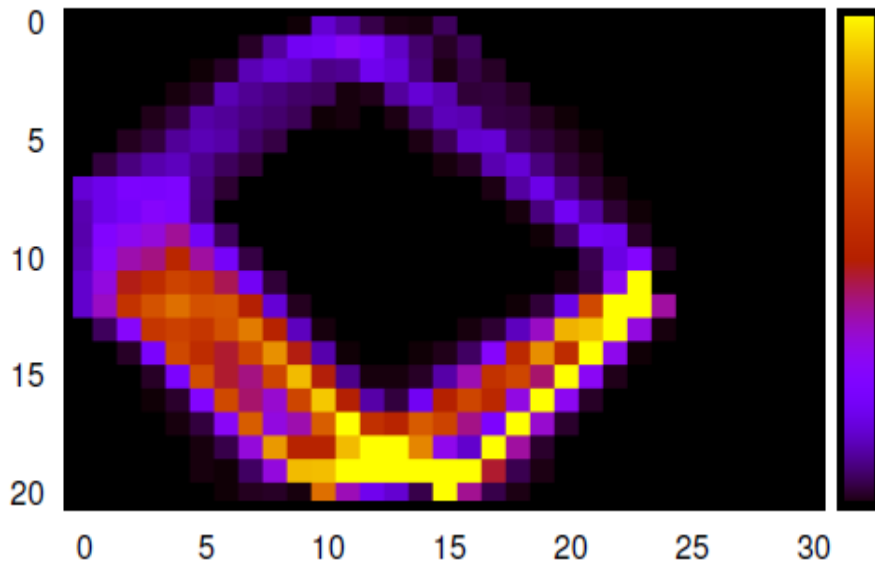


High Variation Fuel World - First 20 Ep. - 15 Trees: $b = 0$

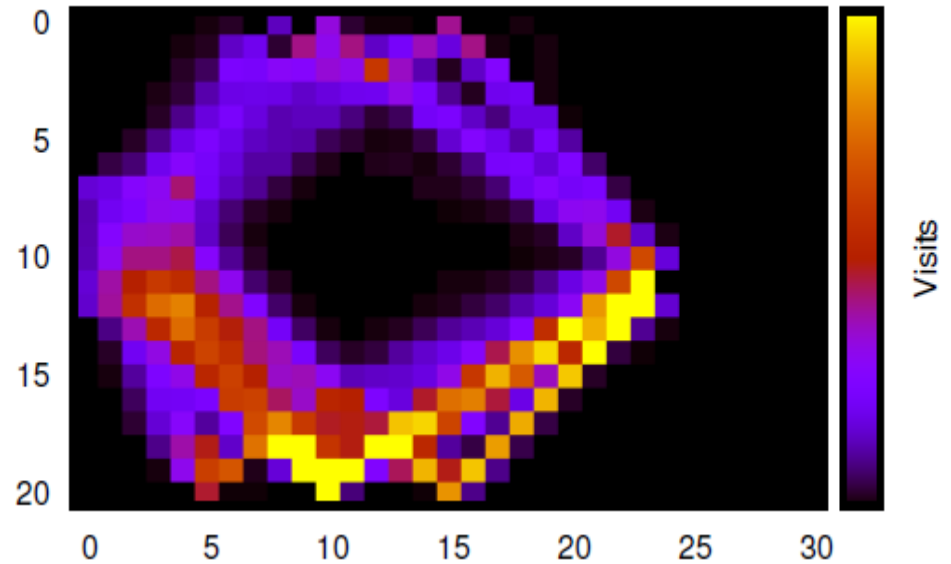


Experiments - Results

Low Variation Fuel World - Final 50 Ep. - 15 Trees: $b = 0$



High Variation Fuel World - Final 50 Ep. - 15 Trees: $b = 0$



The End



Thank you for your Attention

Discussion

Any questions ?