

Bandits, Global Optimization, Active Learning, and Bayesian RL – understanding the common ground

Marc Toussaint

Machine Learning & Robotics Lab – University of Stuttgart
marc.toussaint@informatik.uni-stuttgart.de

Autonomous Learning Summer School, Leipzig, Sep 2014

- This does not focus on own work! It's really a lecture...

The goal is to understand sequential decision problems in which decisions equally influence the learning progress as well as rewards/states.

- Bandits, Global Optimization, Active Learning, and Bayesian RL are instances of this. The perspective taken is simple: All of these problems are eventually Markovian processes in belief space
- For instance, you'll learn what 'optimal optimization' is

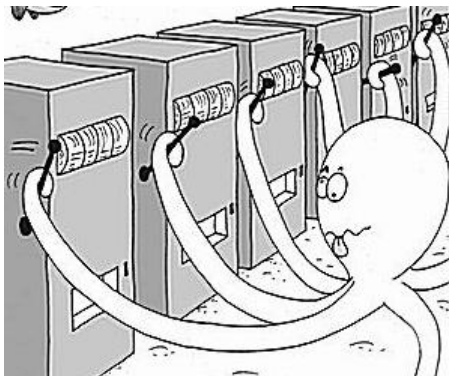
Disclaimer: Whenever I say "optimal" I mean "Bayes optimal" (we always assume having priors $P(\theta)$)

Outline

- Problems covered:
 - Bandits
 - Global optimization
 - Active learning
 - Bayesian RL
 - Monte Carlo Tree Search (MTCS)
- Methods covered (interweaved with the above):
 - Belief planning
 - Upper Confidence Bound (UCB)
 - Expected Improvement, probability of improvement
 - Predictive Entropy, Uncertainty Sampling, Shannon Information
 - Bayesian exploration bonus, R_{\max}
 - Monte Carlo Tree Search (MCTS; UCT)

Bandits

Bandits



- There are n machines.
- Each machine i returns a reward $y \sim P(y; \theta_i)$
The machine's parameter θ_i is unknown

Bandits

- Let $a_t \in \{1, \dots, n\}$ be the choice of machine at time t
Let $y_t \in \mathbb{R}$ be the outcome with mean $\langle y_{a_t} \rangle$
- A policy or strategy maps all the history to a new choice:

$$\pi : [(a_1, y_1), (a_2, y_2), \dots, (a_{t-1}, y_{t-1})] \mapsto a_t$$

- Problem: Find a policy π that

$$\max \left\langle \sum_{t=1}^T y_t \right\rangle$$

or

$$\max \langle y_T \rangle$$

or other objectives like discounted infinite horizon $\max \langle \sum_{t=1}^{\infty} \gamma^t y_t \rangle$

Exploration, Exploitation

- “Two effects” of choosing a machine:
 - You collect more data about the machine \rightarrow knowledge
 - You collect reward
- For example
 - Exploration: Choose the next action a_t to $\min \langle H(b_t) \rangle$
 - Exploitation: Choose the next action a_t to $\max \langle y_t \rangle$

The Belief State

- “Knowledge” can be represented in two ways:
 - as the full history

$$h_t = [(a_1, y_1), (a_2, y_2), \dots, (a_{t-1}, y_{t-1})]$$

- as the **belief**

$$b_t(\theta) = P(\theta|h_t)$$

where θ are the unknown parameters $\theta = (\theta_1, \dots, \theta_n)$ of all machines

- In the bandit case:
 - The belief factorizes $b_t(\theta) = P(\theta|h_t) = \prod_i b_t(\theta_i|h_t)$

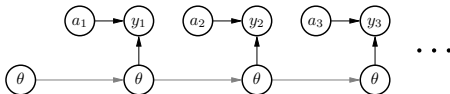
e.g. for binary bandits, $\theta_i = p_i$, with prior $\text{Beta}(p_i|\alpha, \beta)$:

$$b_t(p_i|h_t) = \text{Beta}(p_i|\alpha + a_{i,t}, \beta + b_{i,t})$$

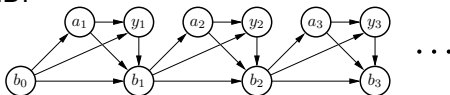
$$a_{i,t} = \sum_{s=1}^{t-1} [a_s = i][y_s = 0], \quad b_{i,t} = \sum_{s=1}^{t-1} [a_s = i][y_s = 1]$$

The Belief MDP

- The process can be modelled as



or as Belief MDP



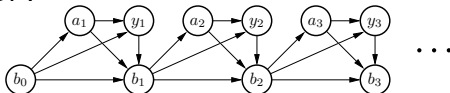
$$P(b'|y, a, b) = \begin{cases} 1 & \text{if } b' = b'_{[b, a, y]} \\ 0 & \text{otherwise} \end{cases}, \quad P(y|a, b) = \int_{\theta_a} b(\theta_a) P(y|\theta_a)$$

- The Belief MDP describes a *different* process: the interaction between the information available to the agent (b_t or h_t) and its actions, where *the agent uses his current belief to anticipate outcomes*, $P(y|a, b)$.

Optimality in the Belief MDP \Rightarrow optimality in the original problem

Optimal policies via Dynamic Programming in Belief Space

- The Belief MDP:



$$P(b'|y, a, b) = \begin{cases} 1 & \text{if } b' = b'_{[b, a, y]} \\ 0 & \text{otherwise} \end{cases}, \quad P(y|a, b) = \int_{\theta_a} b(\theta_a) P(y|\theta_a)$$

- Belief Planning: Dynamic Programming on the value function

$$\begin{aligned} \forall_b : V_{t-1}(b) &= \max_{\pi} \left\langle \sum_{t=t}^T y_t \right\rangle \\ &= \max_{\pi} \left[\langle y_t \rangle + \left\langle \sum_{t=t+1}^T y_t \right\rangle \right] \\ &= \max_{a_t} \int_{y_t} P(y_t|a_t, b) \left[y_t + V_t(b'_{[b, a_t, y_t]}) \right] \end{aligned}$$

$$V_t^*(h) := \operatorname{argmax}_{\pi} \int_{\theta} P(\theta|h) V_t^{\pi, \theta}(h) \quad (1)$$

$$= \operatorname{argmax}_{\pi} \int_{\theta} P(\theta|h) \max_a \left[R(a, h) + \int_{h'} P(h'|h, a, \theta) V_{t+1}^{\pi, \theta}(h') \right] \quad (2)$$

$$V_t^*(b) = \operatorname{argmax}_{\pi} \int_{\theta} b(\theta) \max_a \left[R(a, b) + \int_{b'} P(b'|b, a, \theta) V_{t+1}^{\pi, \theta}(b') \right] \quad (3)$$

$$= \operatorname{argmax}_{\pi} \max_a \int_{\theta} \int_{b'} b(\theta) P(b'|b, a, \theta) \left[R(a, b) + V_{t+1}^{\pi, \theta}(b') \right] \quad (4)$$

$$P(b'|b, a, \theta) = \int_y P(b', y|b, a, \theta) \quad (5)$$

$$= \int_y \frac{P(\theta|b, a, b', y) P(b', y|b, a)}{P(\theta|b, a)} \quad (6)$$

$$= \int_y \frac{b'(\theta) P(b', y|b, a)}{b(\theta)} \quad (7)$$

$$V_t^*(b) = \operatorname{argmax}_{\pi} \max_a \int_{\theta} \int_{b'} \int_y b(\theta) \frac{b'(\theta) P(b', y|b, a)}{b(\theta)} \left[R(a, b) + V_{t+1}^{\pi, \theta}(b') \right] \quad (8)$$

$$= \operatorname{argmax}_{\pi} \max_a \int_{b'} \int_y P(b', y|b, a) \left[R(a, b) + \int_{\theta} b'(\theta) V_{t+1}^{\pi, \theta}(b') \right] \quad (9)$$

$$= \operatorname{argmax}_{\pi} \max_a \int_y P(y|b, a) \left[R(a, b) + \int_{\theta} b'_{[b, a, y]}(\theta) V_{t+1}^{\pi, \theta}(b'_{[b, a, y]}) \right] \quad (10)$$

$$= \max_a \int_y P(y|b, a) \left[R(a, b) + V_{t+1}^*(b'_{[b, a, y]}) \right] \quad (11)$$

Optimal policies

- The value function assigns a value (maximal achievable expected return) to a state of knowledge
- Optimal policies “navigate through belief space”
 - This automatically implies/combines “exploration” and “exploitation”
 - There is no need to explicitly address “exploration vs. exploitation” or decide for one against the other. Optimal policies will automatically do this.
- The optimal policy is greedy w.r.t. the value function (in the sense of the \max_{a_t} above)
- Computationally heavy: b_t is a probability distribution, V_t a function over probability distributions
- The term $\int_{y_t} P(y_t|a_t, b_{t-1}) \left[y_t + V_t(b_{t-1}[a_t, y_t]) \right]$ is related to the *Gittins Index*: it can be computed for each bandit separately.

Exercise

- Consider 3 binary bandits for $T = 10$.
 - How “large” is the belief space? What numbers do you need to store a belief?

Exercise

- Consider 3 binary bandits for $T = 10$.
 - How “large” is the belief space? What numbers do you need to store a belief?
The belief is 3 Beta distributions $\text{Beta}(p_i | \alpha + a_i, \beta + b_i) \rightarrow 6$ integers
 $T = 10 \rightarrow$ each integer ≤ 10

Exercise

- Consider 3 binary bandits for $T = 10$.
 - How “large” is the belief space? What numbers do you need to store a belief?
The belief is 3 Beta distributions $\text{Beta}(p_i | \alpha + a_i, \beta + b_i) \rightarrow 6$ integers
 $T = 10 \rightarrow$ each integer ≤ 10
 - How “large” is the value function $V_t(b_t)$? How many numbers to store $V_t(b_t)$?

Exercise

- Consider 3 binary bandits for $T = 10$.
 - How “large” is the belief space? What numbers do you need to store a belief?
The belief is 3 Beta distributions $\text{Beta}(p_i | \alpha + a_i, \beta + b_i) \rightarrow 6$ integers
 $T = 10 \rightarrow$ each integer ≤ 10
 - How “large” is the value function $V_t(b_t)$? How many numbers to store $V_t(b_t)$?
 $V_t(b_t)$ is a function over $\{0, \dots, 10\}^6 \rightarrow 1$ Mio. numbers $\in \mathbb{R}$

Exercise

- Consider 3 binary bandits for $T = 10$.
 - How “large” is the belief space? What numbers do you need to store a belief?
The belief is 3 Beta distributions $\text{Beta}(p_i | \alpha + a_i, \beta + b_i) \rightarrow 6$ integers
 $T = 10 \rightarrow$ each integer ≤ 10
 - How “large” is the value function $V_t(b_t)$? How many numbers to store $V_t(b_t)$?
 $V_t(b_t)$ is a function over $\{0, \dots, 10\}^6 \rightarrow 1$ Mio. numbers $\in \mathbb{R}$
Many states cannot be visited (integers need to sum up)
Only very few transitions are possible (incrementing integers)

Greedy heuristic: Upper Confidence Bound (UCB)

-
- 1: Initialization: Play each machine once
 - 2: **repeat**
 - 3: Play the machine i that maximizes $\hat{y}_i + \beta \sqrt{\frac{2 \ln n}{n_i}}$
 - 4: **until**
-

\hat{y}_i is the average reward of machine i so far

n_i is how often machine i has been played so far

$n = \sum_i n_i$ is the number of rounds so far

β is often chosen as $\beta = 1$

See *Finite-time analysis of the multiarmed bandit problem*, Auer, Cesa-Bianchi & Fischer, Machine learning, 2002.

UCB algorithms

- UCB algorithms determine a **confidence interval** such that

$$\hat{y}_i - \sigma_i < \langle y_i \rangle < \hat{y}_i + \sigma_i$$

with high probability.

UCB chooses the upper bound of this confidence interval

- *Optimism in the face of uncertainty*
- Strong bounds on the regret (sub-optimality) of UCB (e.g. Auer et al.)

Exercise

- Data so far:
Machine A: 8, 7, 12, 13, 11, 9
Machine B: 8, 12
Machine C: 5, 13

Which one do you choose next?

Exercise

- Data so far:

Machine A: 8, 7, 12, 13, 11, 9

Machine B: 8, 12

Machine C: 5, 13

Which one do you choose next?

Machine A: $10 \pm 2.16/\sqrt{6}$

Machine B: $10 \pm 2/\sqrt{2}$

Machine C: $9 \pm 4/\sqrt{2}$

Conclusions

- The bandit problem is an archetype for
 - Sequential decision making
 - Decisions that influence knowledge as well as rewards/states
 - Exploration/exploitation
- The same aspects are inherent also in global optimization, active learning & RL
- Belief Planning in principle gives the optimal solution
- Greedy Heuristics (UCB) are computationally much more efficient and guarantee bounded regret. MCTS is also applicable

Further reading

- ICML 2011 Tutorial *Introduction to Bandits: Algorithms and Theory*, Jean-Yves Audibert, Rémi Munos
- *Finite-time analysis of the multiarmed bandit problem*, Auer, Cesa-Bianchi & Fischer, Machine learning, 2002.
- *On the Gittins Index for Multiarmed Bandits*, Richard Weber, Annals of Applied Probability, 1992.
Optimal Value function is submodular.

Global Optimization

Global Optimization

- Let $x \in \mathbb{R}^n$, $f : \mathbb{R}^n \rightarrow \mathbb{R}$, find

$$\min_x f(x)$$

(I neglect constraints $g(x) \leq 0$ and $h(x) = 0$ here – but could be included.)

- Blackbox optimization: find optimum by sampling values $y_t = f(x_t)$
No access to ∇f or $\nabla^2 f$
Observations may be noisy $y \sim \mathcal{N}(y \mid f(x_t), \sigma)$

Global Optimization = infinite bandits

- In global optimization $f(x)$ defines a reward for every $x \in \mathbb{R}^n$
 - Instead of a finite number of actions a_t we now have x_t
- Optimal Optimization could be defined as: find $\pi : h_t \mapsto x_t$ that

$$\min \left\langle \sum_{t=1}^T f(x_t) \right\rangle$$

or

$$\min \langle f(x_T) \rangle$$

Gaussian Processes as belief

- The unknown “world property” is the function $\theta = f$
- Given a Gaussian Process prior $GP(f|\mu, C)$ over f and a history

$$D_t = [(x_1, y_1), (x_2, y_2), \dots, (x_{t-1}, y_{t-1})]$$

the belief is

$$b_t(f) = P(f | D_t) = \text{GP}(f|D_t, \mu, C)$$

$$\text{Mean}(f(x)) = \hat{f}(x) = \kappa(x)(\mathbf{K} + \sigma^2\mathbf{I})^{-1}\mathbf{y} \quad \text{response surface}$$

$$\text{Var}(f(x)) = \hat{\sigma}(x) = k(x, x) - \kappa(x)(\mathbf{K} + \sigma^2\mathbf{I}_n)^{-1}\kappa(x) \quad \text{confidence interval}$$

- Side notes:
 - Don't forget that $\text{Var}(y^*|x^*, D) = \sigma^2 + \text{Var}(f(x^*)|D)$
 - We can also handle discrete-valued functions f using GP classification

Optimal optimization via belief planning

- As for bandits it holds

$$\begin{aligned} V_{t-1}(b_{t-1}) &= \max_{\pi} \left\langle \sum_{t=t}^T y_t \right\rangle \\ &= \max_{x_t} \int_{y_t} P(y_t | x_t, b_{t-1}) \left[y_t + V_t(b_{t-1}[x_t, y_t]) \right] \end{aligned}$$

$V_{t-1}(b_{t-1})$ is a function over the GP-belief!

If we could compute $V_{t-1}(b_{t-1})$ we “optimally optimize”

- I don't know of a minimalistic case where this might be feasible

Greedy 1-step heuristics

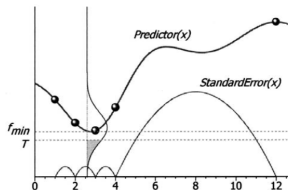


Figure 14. Using kriging, we can estimate the probability that sampling at a given point will 'improve' our solution, in the sense of yielding a value that is equal or better than some target T .

from Jones (2001)

- Maximize Probability of Improvement (MPI)

$$x_t = \operatorname{argmax}_x \int_{-\infty}^{y^*} \mathcal{N}(y|\hat{f}(x), \hat{\sigma}(x))$$

- Maximize Expected Improvement (EI)

$$x_t = \operatorname{argmax}_x \int_{-\infty}^{y^*} \mathcal{N}(y|\hat{f}(x), \hat{\sigma}(x)) (y^* - y)$$

- Maximize UCB

$$x_t = \operatorname{argmin}_x \hat{f}(x) - \beta_t \hat{\sigma}(x)$$

(Often, $\beta_t = 1$ is chosen. UCB theory allows for better choices. See Srinivas et al. citation below.)

From Srinivas et al., 2012:

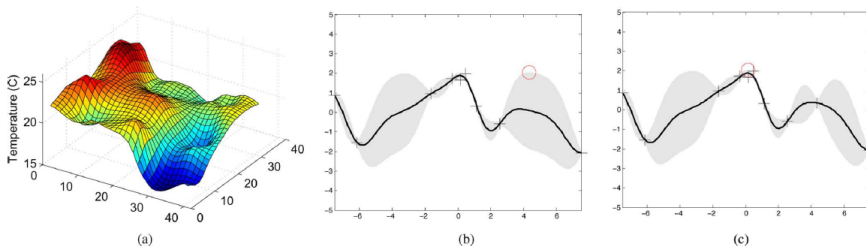


Fig. 2. (a) Example of temperature data collected by a network of 46 sensors at Intel Research Berkeley. (b) and (c) Two iterations of the GP-UCB algorithm. The dark curve indicates the current posterior mean, while the gray bands represent the upper and lower confidence bounds which contain the function with high probability. The “+” mark indicates points that have been sampled before, while the “o” mark shows the point chosen by the GP-UCB algorithm to sample next. It samples points that are either (b) uncertain or have (c) high posterior mean.

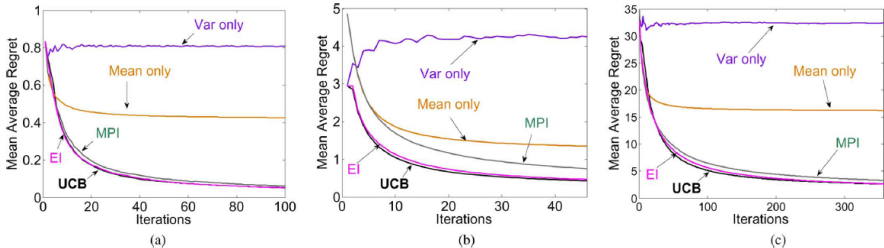


Fig. 6. Mean average regret: GP-UCB and various heuristics on (a) synthetic and (b, c) sensor network data.

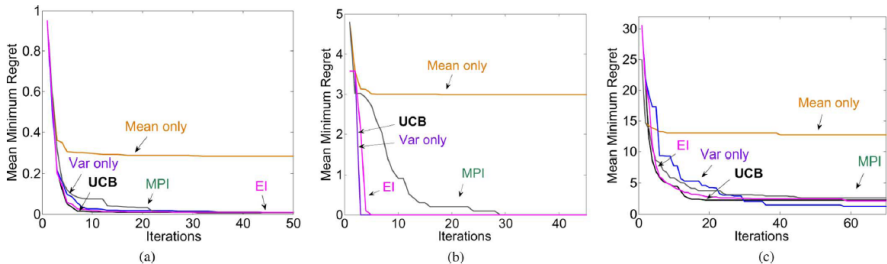


Fig. 7. Mean minimum regret: GP-UCB and various heuristics on (a) synthetic, and (b, c) sensor network data.

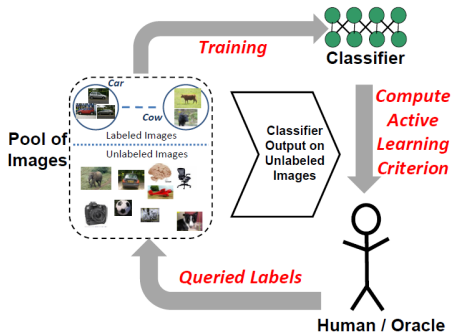
Further reading

- Classically, such methods are known as *Kriging*
- *Information-theoretic regret bounds for gaussian process optimization in the bandit setting* Srinivas, Krause, Kakade & Seeger, Information Theory, 2012.
- *Efficient global optimization of expensive black-box functions.* Jones, Schonlau, & Welch, Journal of Global Optimization, 1998.
- *A taxonomy of global optimization methods based on response surfaces* Jones, Journal of Global Optimization, 2001.
- *Explicit local models: Towards optimal optimization algorithms*, Poland, Technical Report No. IDSIA-09-04, 2004.

Active Learning

Example

Active learning with gaussian processes for object categorization.
Kapoor, Grauman, Urtasun & Darrell, ICCV 2007.



Active Learning

- In standard ML, a data set $D_t = \{(x_s, y_s)\}_{s=1}^{t-1}$ is given.
In active learning, the learning agent sequentially decides on each x_t
– where to collect data
- Generally, the aim of the learner should be to learn as fast as possible, e.g. minimize predictive error
- Finite horizon T predictive error problem:
Given $P(x^*)$, find a policy $\pi : D_t \mapsto x_t$ that

$$\min \langle -\log P(y^*|x^*, D_T) \rangle_{y^*, x^*, D_T; \pi}$$

This also can be expressed as *predictive entropy*:

$$\begin{aligned} \langle -\log P(y^*|x^*, D_T) \rangle_{y^*, x^*} &= \left\langle -\int_{y^*} P(y^*|x^*, D_T) \log P(y^*|x^*, D_T) \right\rangle_{x^*} \\ &= \langle H(y^*|x^*, D_T) \rangle_{x^*} =: H(f|D_T) \end{aligned}$$

- Find a policy that $\min \langle H(f|D_T) \rangle_{D_T; \pi}$

Gaussian Processes as belief

- Again, the unknown “world property” is the function $\theta = f$
- We can use a Gaussian Process to represent the belief

$$b_t(f) = P(f \mid D_t) = \text{GP}(f \mid D_t, \mu, C)$$

Optimal Active Learning via belief planning

- The only difference to global optimization is the reward.
In active learning it is the predictive entropy: $-H(f|D_T)$
- Dynamic Programming:

$$V_T(b_T) = -H(b_T) , \quad H(b) := \langle H(y^*|x^*, b) \rangle_{x^*}$$
$$V_{t-1}(b_{t-1}) = \max_{x_t} \int_{y_t} P(y_t|x_t, b_{t-1}) V_t(b_{t-1}[x_t, y_t])$$

- Computationally intractable

Greedy 1-step heuristic

- The simplest greedy policy is 1-step Dynamic Programming:
Directly maximize immediate expected reward, i.e., minimizes $H(b_{t+1})$.

$$\pi : b_t(f) \mapsto \operatorname{argmin}_{x_t} \int_{y_t} P(y_t|x_t, b_t) H(b_t[x_t, y_t])$$

- For GPs, you reduce the entropy most if you choose x_t where the current predictive variance is highest:

$$\operatorname{Var}(f(x)) = k(x, x) - \kappa(x)(\mathbf{K} + \sigma^2 \mathbf{I}_n)^{-1} \kappa(x)$$

This is referred to as *uncertainty sampling*

- Note, if we fix hyperparameters:
 - This variance is *independent* of the observations y_t , only the set D_t matters!
 - The order of data points also does not matter
 - You can pre-optimize a set of “grid-points” for the kernel – and play them in any order

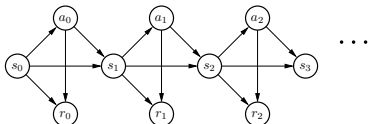
Further reading

- *Active learning literature survey*. Settles, Computer Sciences Technical Report 1648, University of Wisconsin-Madison, 2009.
- *Bayesian experimental design: A review*. Chaloner & Verdinelli, Statistical Science, 1995.
- *Active learning with statistical models*. Cohn, Ghahramani & Jordan, JAIR 1996.
- ICML 2009 Tutorial on *Active Learning*, Sanjoy Dasgupta and John Langford http://hunch.net/~active_learning/

Bayesian Reinforcement Learning

Markov Decision Process

- Other than the previous cases, actions now influence a world state



- initial state distribution $P(s_0)$
 - transition probabilities $P(s'|s, a)$
 - reward probabilities $P(r|s, a)$
 - agent's policy $P(a|s; \pi)$
- Planning in MDPs: Given knowledge of $P(s'|s, a)$, $P(r|s, a)$ and $P(y|s, a)$, find a policy $\pi : s_t \mapsto a_t$ that maximizes the discounted infinite horizon return $\langle \sum_{t=0}^{\infty} \gamma^t r_t \rangle$:

$$V(s) = \max_a \left[\mathbf{E}(r|s, a) + \gamma \sum_{s'} P(s' | s, a) V(s') \right]$$

Bayesian RL: The belief state

- In *Reinforcement Learning* we do not know the world

Unknown MDP parameters $\theta = (\theta_s, \theta_{s'sa}, \theta_{rsa})$

(for $P(s_0), P(s'|s, a), P(r|s, a)$)

- “Knowledge” can be represented in two ways:
 - as the full history

$$h_t = [(s_0, a_0, r_0), \dots, (s_{t-1}, a_{t-1}, r_{t-1}), (s_t)]$$

- as the belief

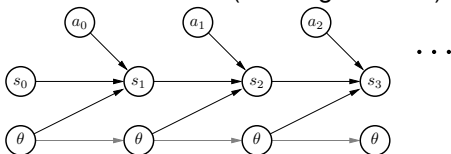
$$b_t(\theta) = P(\theta|h_t)$$

where θ are all the unknown parameters

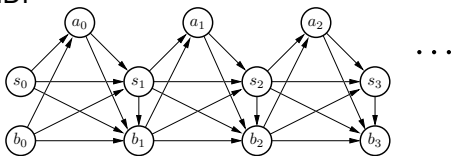
- In the case of discrete MDPs
 - θ are CPTs (conditional probability tables)
 - Assuming Dirichlet priors over CPTs, the exact posterior is a Dirichlet
 - Amounts to counting transitions

Optimal policies

- The process can be modelled as (omitting rewards)



or as Belief MDP



$$P(b'|s', s, a, b) = \begin{cases} 1 & \text{if } b' = b[s', s, a] \\ 0 & \text{otherwise} \end{cases}, \quad P(s'|s, a, b) = \int_{\theta} b(\theta) P(s'|s, a, \theta)$$

$$V(b, s) = \max_a \left[\mathbf{E}(r|s, a, b) + \sum_{s'} P(s'|a, s, b) V(s', b') \right]$$

- Dynamic programming can be approximated (Poupart et al.)

Heuristics

- As with UCB, choose estimators for R^* , P^* that are optimistic/over-confident

$$V_t(s) = \max_a \left[R^* + \sum_{s'} P^*(s'|s, a) V_{t+1}(s') \right]$$

- Rmax:

$$- R^*(s, a) = \begin{cases} R_{\max} & \text{if } \#_{s,a} < n \\ \hat{\theta}_{rsa} & \text{otherwise} \end{cases}, \quad P^*(s'|s, a) = \begin{cases} \delta_{s's^*} & \text{if } \#_{s,a} < n \\ \hat{\theta}_{s'sa} & \text{otherwise} \end{cases}$$

- Guarantees over-estimation of values, polynomial PAC results!
- Read about “KWIK-Rmax”! (Li, Littman, Walsh, Strehl, 2011)
- Bayesian Exploration Bonus (BEB), Kolter & Ng (ICML 2009)
 - Choose $P^*(s'|s, a) = P(s'|s, a, b)$ integrating over the current belief $b(\theta)$ (non-over-confident)
 - But choose $R^*(s, a) = \hat{\theta}_{rsa} + \frac{\beta}{1 + \alpha_0(s, a)}$ with a hyperparameter $\alpha_0(s, a)$, over-estimating return
- Confidence intervals for V -/ Q -function (Kealbling '93, Dearden et al. '99)

Further reading

- ICML-07 Tutorial on Bayesian Methods for Reinforcement Learning
<https://cs.uwaterloo.ca/~ppoupart/ICML-07-tutorial-Bayes-RL.html>
Esp. part 3: Model-based Bayesian RL (Pascal Poupart); and the methods cited on slide 22
- *Optimal learning: Computational procedures for Bayes-adaptive Markov decision processes*. Duff, Doctoral dissertation, University of Massachusetts Amherst, 2002.
- *An analytic solution to discrete Bayesian reinforcement learning*. Poupart, Vlassis, Hoey, & Regan (ICML 2006)
- KWIK-Rmax: *Knows what it knows: a framework for self-aware learning*. Li, Littman, Walsh & Strehl, Machine learning, 2011.
- Bayesian Exploration Bonus: *Near-Bayesian exploration in polynomial time*. Kolter & Ng, ICML 2009.
- The “interval exploration method” described in *Reinforcement learning: A survey*. Kaelbling, Littman & Moore, arXiv preprint cs/9605103, 1996.

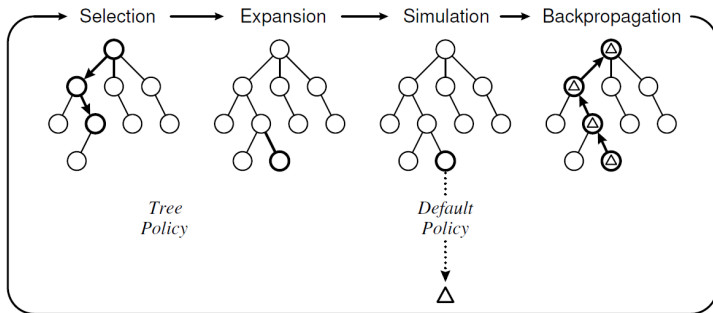
Monte Carlo Tree Search (MCTS)

Monte Carlo Tree Search (MCTS)

- MCTS triggered a little revolution...
- MCTS is very successful on Computer Go and other games
- MCTS is rather simple to implement
- MCTS is very general: applicable on any discrete domain

- Key paper:
Kocsis & Szepesvári: *Bandit based Monte-Carlo Planning*, ECML 2006.
- Survey paper:
Browne et al.: *A Survey of Monte Carlo Tree Search Methods*, 2012.
- POMDPs:
Silver & Veness: *Monte-Carlo Planning in Large POMDPs*, NIPS 2010
- Tutorial presentation:
<http://web.engr.oregonstate.edu/~afern/icaps10-MCP-tutorial.ppt>

Basic MCTS scheme



from Browne et al.

-
- 1: start tree $V = \{v_0\}$
 - 2: **while** within computational budget **do**
 - 3: $v_l \leftarrow \text{TREEPOLICY}(V)$ chooses a leaf of V
 - 4: append v_l to V
 - 5: $\Delta \leftarrow \text{ROLLOUTPOLICY}(V)$ rolls out a full simulation, with return Δ
 - 6: $\text{BACKUP}(v_l, \Delta)$ updates the values of all parents of v_l
 - 7: **end while**
 - 8: return best child of v_0
-

Growing the tree as a sequential decision problem

- We talk here about the internal planning process!
- Deciding to allocate resources to grow the tree in a certain direction (the `TREEPOLICY`) is a decision!
Growing the full tree a sequential decision problem
- What would be the optimal way to make growing decisions?
→ A problem of planning within the planning algorithm...
- The optimal solution is of course infeasible, but...

Upper Confidence Tree (UCT)

- UCT uses UCB to realize the `TREEPOLICY`, i.e. to decide where to expand the tree
- `BACKUP` updates all parents of v_l as
$$n(v) \leftarrow n(v) + 1 \quad (\text{count how often has it been played})$$
$$Q(v) \leftarrow Q(v) + \Delta \quad (\text{sum of rewards received})$$
- `TREEPOLICY` chooses child nodes based on UCB:

$$\operatorname{argmax}_{v' \in \partial(v)} \frac{Q(v')}{n(v')} + \beta \sqrt{\frac{2 \ln n(v)}{n(v')}}}$$

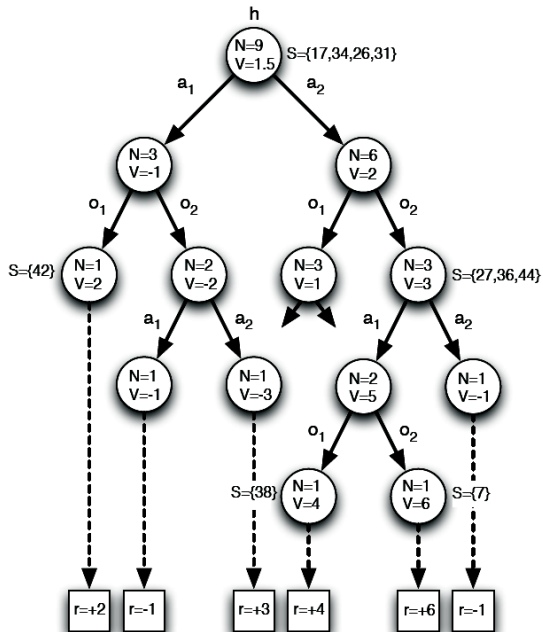
or choose v' if $n(v') = 0$

- In games use a “negamax” backup: While iterating upward, flip sign
$$\Delta \leftarrow -\Delta$$
 in each iteration

Issues when applying MCTS ideas to POMDPs

- key paper:
Silver & Veness: *Monte-Carlo Planning in Large POMDPs*, NIPS 2010
- MCTS is based on generating rollouts using a simulator
 - Rollouts need to start at a specific *state* s_t
 - Nodes in our tree need to have states associated, to start rollouts from
- At any point in time, the agent has only the history $h_t = (y_{0:t}, a_{0:t-1})$ to decide on an action
 - The agent wants to estimate the Q-funcion $Q(h_t, a_t)$
 - Nodes in our tree need to have a history associated
- Nodes in the search tree will
 - maintain $n(v)$ and $Q(v)$ as before
 - have a history $h(v)$ attached
 - have a *set* of states $\mathcal{S}(v)$ attached

MCTS applied to POMDPs



from Silver & Veness

MCTS applied to POMDPs

- For each rollout:
 - Choose a *random* world state $s_0 \sim \mathcal{S}(v_0)$ from the set of states associated to the root v_0 ; initialize the simulator with this s_0
 - Use a TREEPOLICY to traverse the current tree; during this, update the state sets $\mathcal{S}(v)$ to contain the world state simulated by the simulator
 - Use a ROLLOUTPOLICY to simulate a full rollout
 - Append a new leaf v_l with novel history $h(v_l)$ and a single state $\mathcal{S}(v_l)$ associated

Discussion

3 points to make

Point 1: Common ground

What bandits, global optimization, active learning, Bayesian RL & POMDPs share

- Sequential decisions
- Markovian w.r.t. belief
- Decisions influence the knowledge as well as rewards/states
- Sometimes described as “exploration/exploitation problems”

Point 2: Optimality

- In all cases, belief planning would yield optimal solutions
→ Optimal Optimization, Optimal Active Learning, etc...
- Even if it may be computationally infeasible, it is important to know conceptually
- Optimal policies “navigate through belief space”
 - This automatically implies/combines “exploration” and “exploitation”
 - There is no need to explicitly address “exploration vs. exploitation” or decide for one against the other. Policies that maximize the single objective of future returns will automatically do this.

Point 3: Greedy (1-step) heuristics

- Also the optimal policy is greedy – w.r.t. the value function!
- “Greedy heuristics” replace the value function by something simpler and more direct to compute, typically 1-step criteria
 - UCB
 - Probability of Improvement, Expected Improvement
 - Expected immediate reward, expected predictive entropy
- Typically they reflect *optimism in the face of uncertainty*
- Regret bounds for UCB on bandits and optimization (Auer et al.; Srinivas et al.)
- Theory on submodularity very strongly motivates greedy heuristics
- In RL: Optimism w.r.t. θ , but planning w.r.t. s
 - Bayesian Exploration Bonus (BEB), Rmax, interval exploration method

Thanks

for your attention!