

Myopic Gittins Exploration Heuristic for PILCO

Rowan McAllister

August 5, 2015

Contents

1	Situation	1
2	Task	2
3	Approximate Solution: Finite-Horizon	2
3.1	Optimistic Approximation	2
3.1.1	Computing Solution	3
3.1.2	Using Solution	4
3.2	Pessimistic Approximation	4
3.2.1	Example Values of $-\lambda$ for Various Inputs	5
3.2.2	Derivatives	5
4	Approximate Solution: Infinite-Horizon	6
4.1	Derivatives	6
4.2	Bounds	7

1 Situation

The simulate function generates cumulative-cost *distributions* \mathcal{C}_θ (with derivatives) for any controller parameterisation $\theta \in \Theta$ we input,

$$\mathcal{C}_\theta \sim \mathcal{N}(\mu_\theta, \sigma_\theta^2) \leftarrow \text{simulate}(\theta) \quad (1)$$

Currently we do gradient decent, to select the optimal θ w.r.t. mean information only:

$$\theta^* \leftarrow \arg \min_{\theta} [\mu_\theta] \quad (2)$$

In addition to cumulative-cost mean information μ_θ , we also have:

- cumulative-cost uncertainty information σ_θ^2 ,
- the number of trials left n .

We possibly also know the equivalent ‘observation noise’ that a sample from a rollout would give us (by an ability to estimate the reduction in uncertainty that would result from a future rollout).

2 Task

How can we find a θ that is optimal, not just w.r.t. mean information, but w.r.t. all:

- μ_θ mean information,
- σ_θ^2 uncertainty information,
- n number of trials left information.

In addition we *might* know:

- $\hat{\sigma}_{\theta,t+1}^2$ estimated uncertainty of the next time iteration (default 0), or
- $\sigma_y^2(\theta)$ observation noise as a function of θ (default 0 for all θ).

3 Approximate Solution: Finite-Horizon

The exact solution involves infinite belief lookahead (intractable). We therefore restrict our solution-space to myopic (one-step) lookahead solutions. Myopic lookahead is an approximation to infinite belief lookahead, which only considers what could plausibly be learned in the next time step, and thereafter assumes the agent will learn nothing further. If we can make the assumption that each Gaussian \mathcal{C}_θ is independent from all other $\{\mathcal{C}_{\theta'} : \theta' \in \Theta \neq \theta\}$, then we apply the Gittin’s index solution.

3.1 Optimistic Approximation

We can make the optimistic assumption that our observation noise σ_y^2 is zero. I.e. we are about to learn everything there is to learn about our Gaussian bandit in the next time step. With prior cumulative-cost distribution \mathcal{C}_θ^t , and applying controller parameterised by θ , and observing the rollout data \mathcal{D}^t , the posterior cumulative-cost (\mathcal{C}_θ^{t+1}) will collapse to a delta function:

$$\mathcal{C}_\theta^t \sim \mathcal{N}(\mu_\theta^t, \sigma_\theta^2) \quad (3)$$

$$\mathcal{D}^t \leftarrow \text{applyController}(\theta) \quad (4)$$

$$\mathcal{C}_\theta^{t+1}|\mathcal{D}^t \sim \mathcal{N}(\mu_\theta^{t+1}, 0) \quad (5)$$

This is to say, we assume that after just one rollout with parameter θ , we will know (with certainty) the cumulative-cost of a controller parameterised by θ . This is optimistic, since the system itself has non-negligible process noise and has non-negligible observations noise.

This optimism of ‘learning everything about \mathcal{C}_θ with just one rollout’ is known in the POMDP literature as a strict upper bound on the true value function, usually denoted to as V_{MDP} or Q_{MDP} . The value function in POMDPs is sublinear w.r.t. the belief simplex, whereas V_{MDP} is linear (hence an upper bound). In this case we have a belief over not the physical state (as per usual in POMDPs), but a belief over the latent value function. If indeed the agent can learn everything in one rollout, then the upper bound is equal to the true value function.

3.1.1 Computing Solution

The solution involves:

1. using a myopic look-ahead tree,
2. computing the Gittins index (exact except for finite-horizon approximation) of our look-ahead tree (which uses number of trials n information),
3. using gradient decent w.r.t. mean *and* uncertainty.

Note we are minimising cumulative-costs here (rather than the usual bandit setting of maximising cumulative-rewards). OK, so we treat \mathcal{C}_θ^t as a bandit, and we bring in a fixed Gittins bandit which offers a constant cumulative-cost λ per trial. For illustrative simplicity, let us initially assume our prior distribution \mathcal{C}_θ^t is a standard Gaussian $c \sim \mathcal{N}(0, 1)$. The Gittins index is the cost of the fixed-bandit that makes us ambivalent between choosing the uncertain bandit and the fixed bandit, i.e.:

$$\underbrace{\mathcal{C}_{\text{MDP}}(\theta)}_{\text{cost of trying uncertain bandit}} = \underbrace{n\lambda}_{\text{cost of staying with fixed bandit}} \quad (6)$$

We now solve for λ . Let $\phi(\cdot)$ be the standard normal distribution, $\Phi(\cdot)$ its cumulative distribution function:

$$n\lambda = \mathcal{C}_{\text{MDP}}(\theta) \quad (7)$$

$$= \int_{-\infty}^{\infty} \phi(c) \min(\lambda, c) dc \quad (8)$$

$$= \int_{-\infty}^{\lambda} \phi(c) \min(\lambda, c) dc + \int_{\lambda}^{\infty} \phi(c) \min(\lambda, c) dc \quad (9)$$

$$= \int_{-\infty}^{\lambda} n\phi(c)c dc + \int_{\lambda}^{\infty} \phi(c)(c + \lambda(n-1)) dc \quad (10)$$

$$= (n-1) \cdot \int_{-\infty}^{\lambda} \phi(c)c dc + \lambda(n-1) \cdot \int_{\lambda}^{\infty} \phi(c) dc \quad (11)$$

$$= -(n-1)\phi(\lambda) + \lambda(n-1)(1 - \Phi(\lambda)) \quad (12)$$

$$\lambda = -(n-1)(\phi(\lambda) + \lambda\Phi(\lambda)) \quad (13)$$

Solve for λ using Newton's method:

$$g_i = 0 = (n-1)(\phi(\lambda) + \lambda\Phi(\lambda)) + \lambda \quad (14)$$

$$g'_i = \frac{dg_i}{d\lambda} = (n-1)(-\lambda\phi(\lambda) + \Phi(\lambda) + \lambda\phi(\lambda)) + 1 \quad (15)$$

$$= (n-1)\Phi(\lambda) + 1 \quad (16)$$

$$\lambda_0 \leftarrow 0 \quad (17)$$

$$\lambda_{i+1} \leftarrow \lambda_i - \frac{g_i}{g'_i} \quad (18)$$

This converges within 5-6 iterations.

3.1.2 Using Solution

OK great, we have this λ value, which is a function of the number of trails left n . How do we select θ^* ?

$$\mathcal{C}_\theta \sim \mathcal{N}(\mu_\theta, \sigma_\theta^2) \leftarrow \text{simulate}(\theta) \quad (19)$$

$$\theta^* \leftarrow \arg \min_{\theta} [\mu_\theta - \lambda(n, \sigma_\theta^2)] \quad (20)$$

Note how eq (20) equation differs with our current approach in PILCO, eq (2). Not only does eq (20) utilise mean information μ_θ , as does eq (20), but also uncertainty information σ_θ^2 , and horizon information (number of trails remaining) n (via function $\lambda(n)$). Also note there are no free parameters!

3.2 Pessimistic Approximation

In this section we discuss the pessimistic approximation, which is more general than the optimistic approximation (since we can always revert to the optimistic approximation by setting observation noise to zero in this section).

Let our current prior at time t is $\mathcal{C}_t \sim \mathcal{N}(\mu_t, \sigma_t^2)$ (dropping the inherent θ index to simplify notation). We therefore anticipate, according to our prior, that we will receive an observation distributed as

$$y_t \sim \mathcal{N}(\mathcal{C}_t, \sigma_y^2) \quad (21)$$

Thus, our new posterior will be

$$\mathcal{C}_{t+1} \sim \mathcal{N}(\mu_{t+1}, \sigma_{t+1}^2), \text{ where} \quad (22)$$

$$\sigma_{t+1}^2 = (\sigma_t^{-2} + \sigma_y^{-2})^{-1} \quad (23)$$

$$\mu_{t+1} = \sigma_{t+1}^2 (\sigma_t^{-2} \mu_t + \sigma_y^{-2} y_t) \quad (24)$$

We now make the pessimistic assumption that the next time iteration is the only thing we will ever learn (according to the normal observation noise) about this particular parameterisation θ . I.e. the mean function will be our expected reward payout from \mathcal{C}_θ all the way up until the horizon.

OK so what will the approximate Gittins Index look like now? First let us note:

$$y_t \sim \mathcal{N}(\mathcal{C}_t, \sigma_y^2) \quad (25)$$

$$\sim \mathcal{N}(\mu_t, \sigma_t^2 + \sigma_y^2) \quad (26)$$

$$\mu_{t+1} = \sigma_{t+1}^2 (\sigma_t^{-2} \mu_t + \sigma_y^{-2} y_t) \quad (27)$$

$$\sim \mathcal{N}(\sigma_{t+1}^2 (\sigma_t^{-2} \mu_t + \sigma_y^{-2} y_t), \sigma_{t+1}^4 \sigma_y^{-4} (\sigma_t^2 + \sigma_y^2)) \quad (28)$$

$$\sim \mathcal{N}(\mu_t, \underbrace{\sigma_t^2 \cdot \frac{\sigma_t^2}{(\sigma_t^2 + \sigma_y^2)}}_{s^2}) \quad (29)$$

For simple notation let $\mu = \mu_{t+1}$. Since the Gittins index can be decomposed into a mean and additive function of variance, then for these calculation we centre \mathcal{C}_t , i.e. $\mu_t = 0$ (and will un-centre it later). OK now the approximate

pessimistic Gittins index:

$$n\lambda = \int_{-\infty}^{\infty} \mathcal{N}(\mu; 0, s^2) \min(\lambda, \mu) d\mu \quad (30)$$

$$= \int_{-\infty}^{\lambda} \mathcal{N}(\mu; 0, s^2) \min(\lambda, \mu) d\mu + \int_{\lambda}^{\infty} \mathcal{N}(\mu; 0, s^2) \min(\lambda, \mu) d\mu \quad (31)$$

$$\begin{aligned} &= \int_{-\infty}^{\lambda} y + (n-1)\mathcal{N}(\mu; 0, s^2)\mu d\mu + \int_{\lambda}^{\infty} \mathcal{N}(\mu; 0, s^2)(y + \lambda(n-1)) d\mu \\ &= (n-1) \cdot \int_{-\infty}^{\lambda} \mathcal{N}(\mu; 0, s^2)\mu d\mu + \lambda(n-1) \cdot \int_{\lambda}^{\infty} \mathcal{N}(\mu; 0, s^2) d\mu \\ &= (n-1) \left[-s\phi\left(\frac{\lambda}{s}\right) + \lambda(1 - \Phi\left(\frac{\lambda}{s}\right)) \right] \end{aligned} \quad (32)$$

$$\lambda = -(n-1)(s\phi\left(\frac{\lambda}{s}\right) + \lambda\Phi\left(\frac{\lambda}{s}\right)) \quad (33)$$

Solve for λ using Newton's method:

$$g_i = 0 = (n-1)(s\phi\left(\frac{\lambda}{s}\right) + \lambda\Phi\left(\frac{\lambda}{s}\right)) + \lambda \quad (34)$$

$$g'_i = \frac{dg_i}{d\lambda} = (n-1)\left(-\frac{\lambda}{s}\phi\left(\frac{\lambda}{s}\right) + \Phi\left(\frac{\lambda}{s}\right) + \frac{\lambda}{s}\phi\left(\frac{\lambda}{s}\right)\right) + 1 \quad (35)$$

$$= (n-1)\Phi\left(\frac{\lambda}{s}\right) + 1 \quad (36)$$

$$\lambda_0 \leftarrow 0 \quad (37)$$

$$\lambda_{i+1} \leftarrow \lambda_i - \frac{g_i}{g'_i} \quad (38)$$

3.2.1 Example Values of $-\lambda$ for Various Inputs

Note n is the number of trials remaining (including the current trial we are considering), and σ_y^2/σ_t^2 is the ratio of the observation noise from observing a rollout and the prior variance.

n	σ_y^2/σ_t^2	0	0.1	1	10	∞
1		0	0	0	0	0
2		0.2760	0.2632	0.1952	0.0832	0
3		0.4363	0.4160	0.3085	0.1316	0
4		0.5492	0.5236	0.3883	0.1656	0
5		0.6360	0.6064	0.4497	0.1918	0
6		0.7065	0.6736	0.4996	0.2130	0
7		0.7658	0.7301	0.5415	0.2309	0
8		0.8168	0.7788	0.5776	0.2463	0
9		0.8616	0.8215	0.6092	0.2598	0
10		0.9015	0.8595	0.6374	0.2718	0

3.2.2 Derivatives

How would λ change for small changes in σ_t^2 ? To answer this, we first define λ as a function of s , i.e. $\lambda(s)$ (noting that s is the only variable that is directly a function of σ_t^2), and determine how λ would respond to changes in s such that

the constraint $g = 0$ is still satisfied. I.e.:

$$g = 0 = (n-1)(s\phi(\frac{\lambda}{s}) + \lambda\Phi(\frac{\lambda}{s})) + \lambda \quad (39)$$

$$\therefore \frac{dg}{ds} = 0 = (n-1)(\phi(\frac{\lambda}{s}) + \frac{d\lambda}{ds}\Phi(\frac{\lambda}{s})) + \frac{d\lambda}{ds} \quad (40)$$

$$\therefore \frac{d\lambda}{ds} = -\frac{\phi(\frac{\lambda}{s})}{\Phi(\frac{\lambda}{s}) + 1/(n-1)} \quad (41)$$

Now note:

$$\frac{d\lambda}{d\sigma_t^2} = \frac{d\lambda}{ds} \cdot \frac{ds}{d\sigma_t^2} \quad (42)$$

$$= \frac{d\lambda}{ds} \cdot \left[\frac{\sqrt{\sigma_t^2 + \sigma_y^2} - s/2}{\sigma_t^2 + \sigma_y^2} \right] \quad (43)$$

4 Approximate Solution: Infinite-Horizon

Similar to solution in Sec. 3.2, but with infinite horizon and discounting γ .

$$\frac{\lambda}{1-\gamma} = \int_{-\infty}^{\lambda} \mathcal{N}(\mu; 0, s^2) \min(\lambda, \mu) d\mu + \int_{\lambda}^{\infty} \mathcal{N}(\mu; 0, s^2) \min(\lambda, \mu) d\mu \quad (44)$$

$$= \int_{-\infty}^{\lambda} y + \mathcal{N}(\mu; 0, s^2) \frac{\mu\gamma}{1-\gamma} d\mu + \int_{\lambda}^{\infty} \mathcal{N}(\mu; 0, s^2) (y + \frac{\lambda\gamma}{1-\gamma}) d\mu$$

$$= -\frac{s\gamma}{1-\gamma} \cdot \phi(\frac{\lambda}{s}) + \frac{\lambda\gamma}{1-\gamma} \cdot (1 - \Phi(\frac{\lambda}{s})) \quad (45)$$

$$\lambda = -\frac{\gamma}{1-\gamma} \cdot (s\phi(\frac{\lambda}{s}) + \lambda\Phi(\frac{\lambda}{s})) \quad (46)$$

Solve for λ using Newton's method:

$$g_i = 0 = \gamma \cdot (s\phi(\frac{\lambda}{s}) + \lambda\Phi(\frac{\lambda}{s})) + \lambda(1-\gamma) \quad (47)$$

$$= \gamma \cdot (s\phi(\frac{\lambda}{s}) + \lambda\Phi(\frac{\lambda}{s}) - \lambda) + \lambda \quad (48)$$

$$g'_i = \frac{dg_i}{d\lambda} = \gamma \cdot (-\frac{\lambda}{s}\phi(\frac{\lambda}{s}) + \Phi(\frac{\lambda}{s}) + \frac{\lambda}{s}\phi(\frac{\lambda}{s}) - 1) + 1 \quad (49)$$

$$= \gamma \cdot (\Phi(\frac{\lambda}{s}) - 1) + 1 \quad (50)$$

$$\lambda_0 \leftarrow 0 \quad (51)$$

$$\lambda_{i+1} \leftarrow \lambda_i - \frac{g_i}{g'_i} \quad (52)$$

4.1 Derivatives

Similar to Sec. 3.2.2, we have:

$$g = 0 = \gamma \cdot (s\phi(\frac{\lambda}{s}) + \lambda\Phi(\frac{\lambda}{s})) + \lambda(1-\gamma) \quad (53)$$

$$\therefore \frac{dg}{ds} = 0 = \gamma \cdot (\phi(\frac{\lambda}{s}) + \frac{d\lambda}{ds}\Phi(\frac{\lambda}{s})) + \frac{d\lambda}{ds}(1-\gamma) \quad (54)$$

$$\therefore \frac{d\lambda}{ds} = -\frac{\phi(\frac{\lambda}{s})}{\Phi(\frac{\lambda}{s}) - 1 + 1/\gamma} \quad (55)$$

Now note:

$$\frac{d\lambda}{d\sigma_t^2} = \frac{d\lambda}{ds} \cdot \frac{ds}{d\sigma_t^2} \quad (56)$$

$$= \frac{d\lambda}{ds} \cdot \left[\frac{\sqrt{\sigma_t^2 + \sigma_y^2} - s/2}{\sigma_t^2 + \sigma_y^2} \right] \quad (57)$$

4.2 Bounds

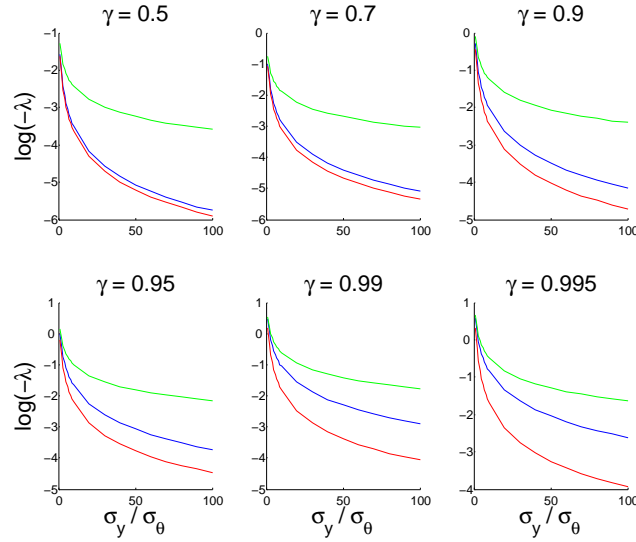


Figure 1: Bounds: The blue line is the exact Gittins index for infinite horizon. The Gittins index is upper bounded by the green line - which is the optimistic approximation to infinite horizon case. The Gittins index is also lower bounded by the red line - which is the pessimistic approximation to infinite horizon case. We set $\sigma_y = 1$.