

**Information Relaxation Bounds for Partially Observed  
Markov Decision Processes**  
*(Online Appendix)*

Martin B. Haugh  
Imperial College Business School  
Imperial College London  
martin.b.haugh@gmail.com

Octavio Ruiz Lacedelli  
Morgan Stanley  
octlacedelli@gmail.com

This is an online appendix including an additional application for the paper "*Information Relaxation Bounds for Partially Observed Markov Decision Processes*" forthcoming in IEEE Transactions on Automatic Control, August 2020.

## G. Application: Personalized Medicine in Mammography Screening

To date there have been relatively few<sup>23</sup> successful medical applications of POMDPs. The reasons for this include the difficulty of determining a suitable objective to optimize, the difficulty of estimating the POMDP parameters and the general difficulty of solving POMDP problems. Recently Ayer et al. [1] proposed a POMDP formulation with the goal of determining an optimal screening policy for breast cancer, the most common cancer among U.S. women according to the American Cancer Society (ACS). The recommendation guidelines provided by the ACS in 2015 [7] is for women with an average risk of breast cancer to take mammograms beginning at age 45, and to continue annually until age 54. Beginning at age 55, they are then recommended to undergo biannual screenings (but they have the opportunity to continue annually if desired) and to continue taking mammograms as long as their life expectancy is at least 10 years. In addition, the ACS indicates that women aged 40 to 44 may choose to begin mammogram screening if desired. In contrast, in 2016 the U.S. Preventive Services Task Force (USPSTF) [9] recommended that women aged 50 to 74 screen biannually using mammography, and they left open the decision for women aged 40-49. In addition, they did not find enough evidence to recommend taking mammograms beyond the age of 75.

In this section we apply the information relaxation approach to the POMDP formulation of Ayer et al. We will use the term decision-maker (DM) to refer to the woman or patient in question but the decision-maker could also refer to a doctor or some other medical professional. We assume the DM has the objective of maximizing her total expected quality-adjusted life years (QALYs). We assume a finite-horizon discrete-time model where the time intervals correspond to six-month periods beginning at age 40 and ending at age 100 so that  $t \in \{0, \dots, 120\}$ . The hidden state space represents the true health state of the patient with  $\mathcal{H} = \{0, 1, 2, 3, 4, 5\}$ . Specifically:

- State 0 represents a cancer-free patient.
- States 1 and 2 indicate the presence of *in situ* and *invasive* cancer, respectively.
- States 3 and 4 represent fully observed absorbing states in which the patient has been diagnosed with *in situ* and *invasive* cancer, respectively, and has begun treatment.
- State 5 is a fully observed absorbing state representing the death of the patient.

Clearly states 3, 4 and 5 can be explicitly observed and are therefore not actually *hidden*. We include them among the set of hidden states, however, to account for the possible transition dynamics of the other hidden states into these absorbing states. The knowledge of being in these hidden absorbing states can then be modeled correctly through noiseless observations of them. We will refer to the subset of hidden states  $\{0, 1, 2\}$  as *pre-cancer* states and the absorbing states  $\{3, 4, 5\}$  as *post-cancer* states.

At each time  $t$ , the DM can choose to either have a mammography screening ( $M$ ) or wait ( $W$ ). If the decision to wait is made, the patient may perform a self-detection screening which will have either a positive or negative result. That is, if through self-detection the patient has reason to be concerned about the presence of cancer, we say the self-test is positive. The possible results of a mammogram are also positive or negative. In the former case, an accurate procedure, e.g. a biopsy, is then prescribed to precisely determine the true cancer status of the patient. If the biopsy result is positive and cancer is found with certainty, the patient will then exit the screening process and move into one of the absorbing states, 3 or 4, to indicate that cancer

---

<sup>23</sup>A review of applications of MDPs and POMDPs to medical decision problems can be found in [8].

treatment has commenced. To code this behavior, Ayer uses hidden state transitions that are functions of the observations. To model this behavior as a conventional POMDP (where hidden state transitions do not depend on observations), we introduce an exit action ( $E$ ) as the only available action after a positive biopsy has been observed. The transition into absorbing state 3 or 4 will now only depend on the current hidden state and the exit action which must be taken if the biopsy result is positive and cancer is found with certainty. The set of possible observations is therefore  $\mathcal{O} = \{R-, R+, B_1, B_2, D\}$  where:

- $R-$  is a negative test result (either from a mammography or self-detection).
- $R+$  is a positive test result (including a negative biopsy if the test was a positive mammogram).
- $B_1$  and  $B_2$  represent in situ cancer and invasive cancer, respectively, and they can be observed via a biopsy following a positive mammogram. If  $B_1$  or  $B_2$  are observed, the action space is then restricted to the exit action  $E$  which transfers the patient to the corresponding absorbing state.
- $D$  represents the death of the patient.

We assume a prior probability distribution,  $\pi_0$ , on the true health-state of the woman at age 40. The transition probabilities of the latent pre-cancer health states are assumed to be age-specific and therefore a function of time  $t$ . We assume that a screening decision does not influence the development of cancer and therefore have  $P_{ij}^t(M) = P_{ij}^t(W)$  for all  $t$  and for all  $i, j \in \mathcal{H}$ . The time  $t$  transition matrices for the screening and wait actions,  $M$  and  $W$ , are then given by

$$P^t(M) = P^t(W) = \begin{bmatrix} p_{00}^t & p_{01}^t & p_{02}^t & 0 & 0 & m_0^t \\ 0 & p_{11}^t & p_{12}^t & 0 & 0 & m_1^t \\ 0 & 0 & p_{22}^t & 0 & 0 & m_2^t \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} \quad (104)$$

where  $m_i^t$  represents the mortality rates for each health-state,  $i$ ,  $p_{01}^t$  and  $p_{02}^t$  represent the in situ and invasive cancer incidence rates, respectively, and  $p_{12}^t$  is the probability that in situ cancer develops into invasive cancer. Recalling that time steps in the POMDP correspond to half-year periods, all rates correspond to effective semiannual rates. Estimates for some of the parameters in (104) were obtained from various sources (see Table 2 below), and we used reasonable assumptions to estimate the parameters for which we could not find external estimates. We note that we have not conducted a full study on the appropriateness<sup>24</sup> of these parameters, but rather we treat them as ballpark estimates in order to illustrate the information relaxation POMDP methodology. Finally, the exit action,  $E$ , will take pre-cancer states to post-cancer treatments with probability 1, i.e.  $P_{1,3}^t(E) = 1$  and  $P_{2,4}^t(E) = 1$ . Since this action is only available to true health-states 1 and 2, we need not define the transitions for other health-states.

<sup>24</sup>Experts in the field of breast cancer could almost certainly provide superior estimates for those parameters where we could not find external estimates.

<sup>25</sup>We approximated the invasive cancer mortality rate by inferring the 6-month mortality rate from the 5 year survival rate (0.897) and used the maximum of this 6-month rate and the average female 6-month mortality for a woman of that age. We assumed that in situ mortality is equal to the female mortality times 1.02 for women of the same age.

<sup>26</sup>The initial risk for an average woman was taken from the breast cancer prevalence rate (0.9462%) and split 80% for invasive cancer and 20% for in situ cancer, as discussed in [11].

Parameter	Source
Mortality $m_0$	SSA Period Life Table, 2013, Female mortality [10]
Mortality $m_1, m_2$	SEER [4] Table 4.13, all stages and all ages <sup>25</sup>
Incidence $p_{01}$	SEER Table 4.12, all races
Incidence $p_{02}$	SEER Table 4.11, all races
Incidence $p_{12}$	Assumed equal to $p_{02}$
Initial risk $\pi_0$	SEER Table 4.24, female 40-49 <sup>26</sup>

Table 2: Sources of the demographic rates for the transition probabilities.

The observation probabilities are determined by the accuracy of the examinations, which are commonly referred to as *specificity* and *sensitivity*. The specificity of a test corresponds to the true negative rate, i.e. the probability that a cancer-free woman obtains a negative test result, while the sensitivity of a test is the true positive rate, i.e. the probability of a positive test result given that the woman has cancer. For each test we employ the age-specific sensitivity and specificity factors that were computed and reported by Ayer et al. They are:

$$\begin{aligned} \text{spec}_t(W) &= 0.92, \quad \forall t & \text{sens}_t(W) &= 0.44, \quad \forall t \\ \text{spec}_t(M) &= \begin{cases} 0.889, & \text{if } t \in \{0, \dots, 19\} \\ 0.893, & \text{if } t \in \{20, \dots, 39\} \\ 0.897, & \text{if } t \geq 40 \end{cases} & \text{sens}_t(M) &= \begin{cases} 0.722, & \text{if } t \in \{0, \dots, 29\} \\ 0.81, & \text{if } t \in \{30, \dots, 59\} \\ 0.862, & \text{if } t \geq 60. \end{cases} \end{aligned}$$

Using these rates, we define the age-specific observation matrices according to

$$B^t(W) = \begin{bmatrix} \text{spec}_t(W) & 1 - \text{spec}_t(W) & 0 & 0 & 0 \\ 1 - \text{sens}_t(W) & \text{sens}_t(W) & 0 & 0 & 0 \\ 1 - \text{sens}_t(W) & \text{sens}_t(W) & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

and

$$B^t(M) = \begin{bmatrix} \text{spec}_t(M) & 1 - \text{spec}_t(M) & 0 & 0 & 0 \\ 1 - \text{sens}_t(M) & 0 & \text{sens}_t(M) & 0 & 0 \\ 1 - \text{sens}_t(M) & 0 & 0 & \text{sens}_t(M) & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

where  $B_{ij}^t(a)$  is the probability of observation  $j \in \mathcal{O}$  when action  $a$  is taken and the hidden state is  $i \in \mathcal{H}$ . Note that the observability of the hidden absorbing states 3, 4 and 5 is made evident through these matrices. It is worth pointing out that once action  $E$  has been chosen, the DM immediately transitions to an absorbing

fully-observable state, and therefore there is no need to define  $B^t(E)$ .

A characteristic of medical decision problems, as pointed out in Ayer et al., is that the observation at time  $t$  is a function of the current action,  $B_{ij}(a) := \mathbb{P}(o_t = j \mid h_t = i, a_t = a)$ , as opposed to a conventional POMDP where the observation is a function of the prior action; see (2). This means that events take place in the following order: given a belief state the DM first takes an action, then immediately observes the result of the action and updates the belief, then a transition takes place and the belief is “carried forward”. This technicality results in a different version of the standard filtering update in which the transition occurs prior to the observation. Nonetheless, filtering in this non-standard form of the POMDP is still<sup>27</sup> a straightforward task. And for the same reason, the natural filtration for the medical decision problem is one where  $\mathcal{F}_t$  is defined to be the  $\sigma$ -algebra generated by  $o_{0:t-1}$ , for  $t \geq 1$ , and with  $\mathcal{F}_0$  defined to be the  $\sigma$ -algebra generated by  $\pi_0$ , the prior distribution on the initial hidden state.

We define the reward obtained at time  $t$ ,  $r_t(h_t, a_t, o_t)$ , as the expected QALYs between times  $t$  and  $t + 1$  that a person in true health-state  $h_t$  would accrue after making decision  $a_t$  and obtaining observation  $o_t$ . Note that although the reward is a function of the as yet unseen observation (see previous paragraph),  $o_t$ , we can instead use<sup>28</sup> its expected value

$$r_t(h_t, a_t) := \mathbb{E}[r_t(h_t, a_t, o_t) \mid h_t] \tag{105}$$

which is easy to calculate and is now in the standard form for a POMDP.

We follow the same calculations as Ayer et al. to define the reward functions. If the patient is in a pre-cancer state  $i \in \{0, 1, 2\}$ , the wait action reward is given by  $r_t(i, W, o_t) = 0.25m_i^t + 0.5(1 - m_i^t)$  where the  $m_i^t$ 's are the (semi-annual) mortality rates given in Table 2. This is the reward for a woman in period  $t$  and true pre-cancer health state,  $i$ , and in fact does not depend on the observation  $o_t$ . Specifically, if death occurs in the next six months (which occurs w.p.  $m_i^t$ ), it is assumed to happen exactly at the three month mark and so the woman will therefore obtain 0.25 years of lifetime. In contrast, if she survives (which occurs w.p.  $1 - m_i^t$ ) she obtains the 0.5 half-years of lifetime in that period.

For the mammography screening action, we subtract a disutility function,  $du(h_t, o_t)$ , from the reward so that  $r_t(h_t, M, o_t) := r_t(h_t, W, o_t) - du(h_t, o_t)$ . The disutility is given a value of 0.5 days for a negative mammogram, two weeks for a true positive mammogram and four weeks for a false positive mammogram. True positive mammograms will in addition force the DM to exit the system in the next period, and provide a lump-sum reward of  $R_t(i) := r_t(i, E)$  for  $i = 1, 2$ . Recall that a true positive mammogram followed by an exit action refers to a woman being accurately diagnosed with cancer and then going into treatment immediately. We expect that a patient under treatment would have a lower remaining expected lifetime than the remaining expected lifetime,  $e_t(0)$  say, of a healthy woman of the same age, but higher than the remaining expected lifetime,  $e_t(i)$  say, of a woman with cancer  $i \in \{1, 2\}$  who is undiagnosed and of the same age. (Note that the expected remaining lifetimes can be calculated using the corresponding mortality rates from times  $t$  to  $T$ .) We therefore assume  $e_t(0) < R_t(i) < e_t(i)$  and in our numerical example, we set  $R_t(i) = 0.5e_t(0) + 0.5e_t(i)$  for  $i = 1, 2$ . We also assume that the absorbing states provide no rewards.

<sup>27</sup>In fact, the timing of observations in this application mirrors that of the multiaccess communication application of Section 8, where an observation occurs immediately after the selection of an action so that each observation is a function of the current hidden state and the *current* action. Therefore, the discussion in Appendix E can be followed to derive the technical details for this section.

<sup>28</sup>We acknowledge a slight abuse of notation here in that we are using the same  $r_t$  to denote time  $t$  rewards  $r_t(h_t, a_t)$ ,  $r_t(h_t, a_t, o_t)$  and  $r_t(\pi_t, a_t)$ . It should be clear from the context what version of the reward we have in mind.

It is perhaps worth noting how the benefit of mammography screening is modelled in our POMDP setting. Specifically, it arises from the possibility of identifying a cancer early and therefore entering treatment and having an expected remaining lifetime that is greater than if the cancer went undiagnosed. The reduced expected lifetime of a woman with an undiagnosed cancer will be reflected via the specific values of the transition and mortality rates of the second and third rows (corresponding to undiagnosed cancer states 1 and 2) in (104). There is a cost to mammography screening, however, which is reflected via the disutility function and so the ultimate goal is to find a policy that trades the benefits of mammography screening off against its disutility.

### G.1. Value Function Approximations

Two methods were used to obtain value function approximations: a QMDP approximation, adapted from the robot navigation problem to include intermediate rewards, and a grid-based approximation. The QMDP approximation is given by

$$\tilde{V}_t(o_{0:t-1}) := \max_{a_t} \sum_{h \in \mathcal{H}} \pi_t(h) V_t^Q(h, a_t) \quad (106)$$

with the understanding that at  $t = 0$ ,  $\tilde{V}_0 := \tilde{V}_0(\pi_0)$ , and where  $V_t^Q$  is the Q-function of the corresponding fully observable MDP formulation, i.e.

$$V_t^Q(h, a) := r_t(h, a) + \sum_{h' \in \mathcal{H}} P_{hh'}(a) V_{t+1}^{\text{MDP}}(h') \quad (107)$$

$$V_t^{\text{MDP}}(h) := \max_{a_t \in \mathcal{A}} V_t^Q(h, a_t) \quad (108)$$

for  $t \in \{0, \dots, T\}$  with terminal condition  $V_{T+1}^{\text{MDP}} := 0$ . Note that the only difference between these definitions and those given for the robot navigation application is the inclusion here of intermediate rewards.

The *grid approximation* corresponds to a point-based value iteration method using a fixed and finite grid approximation of the belief space,  $\Pi$  (see [5, 3]). A standard approximation tool in dynamic programming is to represent an infinite state space as a finite grid of points,  $P \subset \Pi$ , and obtain an AVF by linear interpolation for points not in  $P$ . Specifically, the AVF is obtained by solving a dynamic program with terminal condition  $\tilde{V}_{T+1} = 0$  and Bellman equation

$$\tilde{V}_t(\pi) = \max_{a_t} \left[ r_t(\pi, a_t) + \sum_o \mathbb{P}_{a_t}(o | \pi) \tilde{V}_{t+1}(f(\pi, a_t, o)) \right] \quad (109)$$

for  $t \in \{0, \dots, T\}$ ,  $\pi \in P$  and where  $r_t(\pi, a) := \sum_h r_t(h, a) \pi(h)$ ,  $f(\pi, a, o)$  is the belief update function, and  $\mathbb{P}_a(o | \pi) := \sum_h \mathbb{P}_a(o | h) \pi(h)$ . Note that in general  $f(\pi, a_t, o)$  will not be an element in  $P$  and so we use linear interpolation to evaluate the AVF at those points. To tie in the grid approximation with our application, we take the 3-dimensional subspace corresponding to the pre-cancer states

$$\tilde{\Pi} := \{ \pi \in \Pi \mid \pi = (\pi_0, \pi_1, \pi_2, 0, 0, 0), \pi_0 + \pi_1 + \pi_2 = 1 \}$$

of the 6-dimensional simplex  $\Pi$ . We call  $\tilde{\Pi}$  the pre-cancer belief space simplex<sup>29</sup> and form a finite grid

<sup>29</sup>Although the dimension of the hidden state space is 6, in reality the uncertainty in the process is entirely restricted to the 3 pre-cancer states. We can therefore reduce our analysis to the 3-dimensional pre-cancer belief space simplex.

$P \subset \tilde{\Pi}$ . We then solve the dynamic program (109) for all elements of  $P$  union the elements  $(0, 0, 0, 1, 0, 0)$ ,  $(0, 0, 0, 0, 1, 0)$  and  $(0, 0, 0, 0, 0, 1)$ . For our application, we use a grid  $P$  with elements  $0.05 \times (i_1, i_2, i_3)$  with  $i_1, i_2, i_3$  integer valued and such that they lie on  $\tilde{\Pi}$ , i.e.  $0.05 \times (i_1 + i_2 + i_3) = 1$ .

We can now generate lower bounds on the optimal value function,  $V_0^*(\pi_0)$ , by simulating the policies that are greedy w.r.t. each of the value function approximations. We will compare the performance of these greedy policies to the official policies recommended by ACS and USPSTF.

## G.2. The Uncontrolled Formulation

The action-independent transition and emission matrices are built using different approaches for each AVF. First, using the fact that the QMDP approximation is a supersolution, we can drop the absolute continuity requirement and set the transition matrices,  $Q^t$ , using (63) and, similarly, we set the uncontrolled emission matrices according to

$$E_{ij}^t \equiv B_{ij} \left( \underset{a \in \mathcal{A}}{\operatorname{argmax}} V_t^Q(i, a) \right). \quad (110)$$

In contrast, there is no guarantee that the AVF based on the grid approximation is a supersolution and so we must satisfy the absolute continuity conditions. To achieve this, we add a small positive quantity  $\epsilon = 0.001$  to each  $Q_{ij}^t$  if  $j$  can be reached from  $i$  under some action, and then normalize the probabilities. Similarly, we add  $\epsilon$  to  $B_{ik}^t$  only if  $k$  can be observed from state  $i$  under some action and again we then normalize the probabilities. This approach allows our transition and emission probabilities to satisfy absolute continuity for the PI relaxation. For the BSPI relaxation we would need to make an additional adjustment (as described in Appendix A.2) but the BSPI results were slightly inferior to the PI results (as was the case with the maze application) and so we don't report them in our numerical results.

## G.3. Numerical Results

We consider two different test cases: case 1 represents a woman at age 40 with an average risk of having cancer and therefore an initial distribution over hidden states given by  $\pi_0 = [0.9905, 0.0019, 0.0076, 0, 0, 0]$ . Case 2 represents a woman at age 40 with a high-risk of having cancer; she has an initial distribution of  $\pi_0 = [0.96, 0.02, 0.02, 0, 0, 0]$ . In Figure 5a we display the lower bounds obtained by simulating each of the four policies, namely the policies recommended USPSTF and ACS, as well as the policies that are greedy w.r.t. the QMDP and grid-based AVFs. We note that the latter two policies outperform the official recommendations of USPSTF and ACS, with the best lower bound coming from the grid approximation.

Figure 5b displays the upper bounds obtained with the PI relaxation using penalties constructed from each of the two AVFs. Since the QMDP AVF is a supersolution and therefore also an upper bound we also plot its value in the figure. As a reference, we also display the value of the best lower bound to obtain a visual representation of the duality gap. The duality gap reduction of the best dual bound with respect to the supersolution is 57% in case 1, and 51% in case 2, or equivalently, 19.7 and 29.6 days respectively.

In Ayer et al., the authors were able to solve the POMDP to optimality using Monahan's algorithm [6] with Eagle's reduction [2]. The authors used an Intel Xeon 2.33 GHz processor with 16 GB RAM for their computations, and were able to solve the problem in 55.95 hours. As with our robot navigation application, we used MATLAB Release 2016b, and a MacOS Sierra with 1.3 GHz Intel Core i5 processor with 4 GB RAM. The numerical results in Table 3 display the running times and other statistics for the various Case

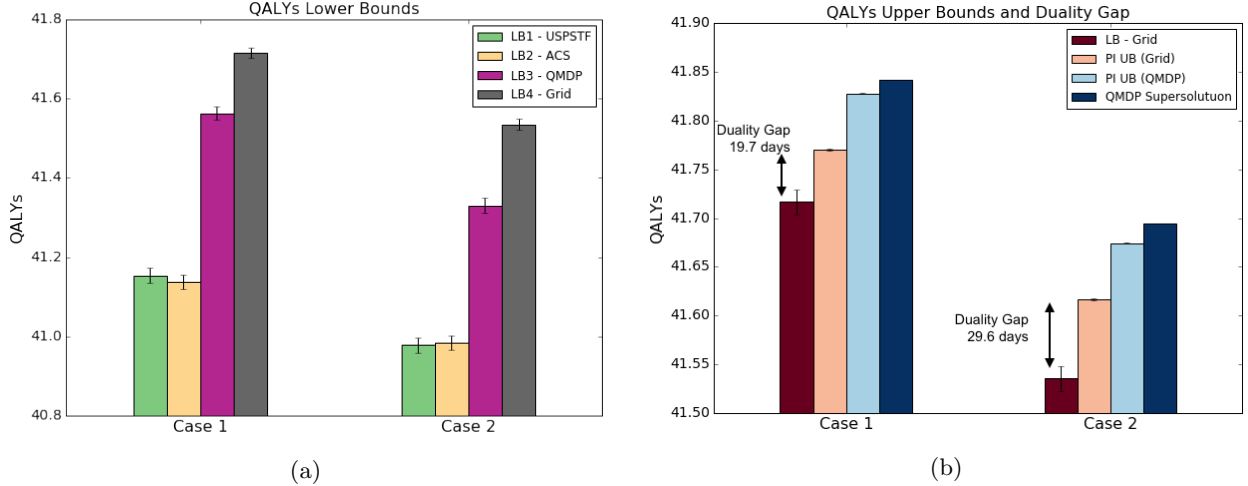


Figure 5: (a) Lower bounds on the optimal value function obtained from simulating the USPSTF and ACS recommended policies as well as policies that are greedy w.r.t. the QMDP and grid-based AVFs. Case 1 corresponds to an average risk 40-year old woman while case 2 corresponds to a high risk 40-year old woman. The vertical lines on each bar represent 95% confidence intervals. (b) Upper bounds on the optimal value function compared to the best lower bound which was obtained by simulating the policy that is greedy w.r.t. the grid-based AVF. The best upper bound was also obtained by constructing penalties for the PI relaxation from the grid-based AVF. The optimal duality gap is displayed in each case.

Bound	Expected value	Standard dev.	Number of paths	Running time
USPSTF (Lower)	41.15	0.0188	400,000	6.75 mins.
ACS (Lower)	41.14	0.0185	400,000	6.79 mins.
Greedy QMDP (Lower)	41.56	0.0160	500	11.1 secs.*
Greedy Grid (Lower)	<b>41.72</b>	0.0128	800,000	35.05 mins
Grid PI (Upper)	<b>41.77</b>	0.0001	100	9.9 secs.
QMDP PI (Upper)	41.83	0.00004	500	8.1 secs.
QMDP Supersolution (Upper)	41.84	-	1	0.02 secs.

\*Lower bound for QMDP greedy strategy was estimated using the penalties as control variates - see Appendix D

Table 3: Summary statistics for the lower and upper bounds for the Case 1 scenario.

1 bounds as well as the best bounds in bold font. As we have noted, our bound approximations result in a very tight duality gap (19.7 days or 0.054 QALYs for an average woman) and we were able to obtain the best lower and upper bounds in 35.05 minutes and 9.9 seconds, respectively, with narrow confidence intervals. So while Ayer et al. were able to solve the problem to optimality, we were able to get provably close to optimality using<sup>30</sup> a slower processor and less RAM with a total runtime that was approximately 2 orders of magnitude smaller. We also note that even tighter bounds information relaxation bounds should be attainable here if so desired using a *partially controlled* formulation as introduced in BH.

<sup>30</sup>We do not know what software Ayer et al. used



## References

- [1] Turgay Ayer, Oguzhan Alagoz, and Natasha K. Stout. Or forum—a pomdp approach to personalize mammography screening decisions. *Operations Research*, 60(5):1019–1034, 2012.
- [2] J N Eagle. The optimal search for a moving target when the search path is constrained. *Operations Research*, 32(5):1107–1115, 1984.
- [3] M. Hauskrecht. Value-function approximations for partially observable markov decision processes. *Journal of Artificial Intelligence Research*, 13:33–94, 2000.
- [4] N. Howlader, A.M. Noone, M. Krapcho, and et al. (editors). SEER Fast Stats, 2009 - 2013. National Cancer Institute. Bethesda, MD., 2016.
- [5] W.S. Lovejoy. A survey of algorithmic methods for partially observed markov decision processes. *Annals of Operations Research*, 28(1):47–65, 1991.
- [6] G E Monahan. A survey of partially observable Markov decision processes: Theory, models, and algorithms. *Management Science*, 28(1):1–16, 1982.
- [7] K.C. Oeffinger, E.H. Fontham, R. Etzioni, and et al. Breast cancer screening for women at average risk: 2015 guideline update from the american cancer society. *JAMA*, 314(15):1599–1614, 2015.
- [8] Andrew J Schaefer, Matthew D Bailey, Steven M Shechter, and Mark S Roberts. Modeling medical treatment using markov decision processes. In *Operations research and health care*, pages 593–612. Springer, 2005.
- [9] A.L. Siu and U.S. preventive services task force. Screening for breast cancer: U.S. preventive services task force recommendation statement. *Annals of Internal Medicine*, 164:279–96, 2016.
- [10] Social Security Administration. Period life table, 2013.
- [11] B. L. Sprague and A. Trentham-Dietz. Prevalence of breast carcinoma in situ in the united states. *JAMA: The Journal of the American Medical Association*, 302(8):846–848, 2009.