# Improving PAC Exploration Using the Median Of Means

**Jason Pazis**
Laboratory for Information and Decision Systems
Massachusetts Institute of Technology
Cambridge, MA 02139, USA
jpazis@mit.edu

**Ronald Parr**
Department of Computer Science
Duke University
Durham, NC 27708
parr@cs.duke.edu

**Jonathan P. How**
Aerospace Controls Laboratory
Department of Aeronautics and Astronautics
Massachusetts Institute of Technology
Cambridge, MA 02139, USA
jhow@mit.edu

## Abstract

We present the first application of the median of means in a PAC exploration algorithm for MDPs. Using the median of means allows us to significantly reduce the dependence of our bounds on the range of values that the value function can take, while introducing a dependence on the (potentially much smaller) variance of the Bellman operator. Additionally, our algorithm is the first algorithm with PAC bounds that can be applied to MDPs with unbounded rewards.

## 1 Introduction

As the reinforcement learning community has shifted its focus from heuristic methods to methods that have performance guarantees, PAC exploration algorithms have received significant attention. Thus far, even the best published PAC exploration bounds are too pessimistic to be useful in practical applications. Even worse, lower bound results [14, 7] indicate that there is little room for improvement.

While these lower bounds prove that there exist pathological examples for which PAC exploration can be prohibitively expensive, they leave the door open for the existence of "well-behaved" classes of problems in which exploration can be performed at a significantly lower cost. The challenge of course is to identify classes of problems that are general enough to include problems of real-world interest, while at the same time restricted enough to have a meaningfully lower cost of exploration than pathological instances.

The approach presented in this paper exploits the fact that while the square of the maximum value that the value function can take ($Q_{\max}^2$) is typically quite large, the variance of the Bellman operator is rather small in many domains of practical interest. For example, this is true in many control tasks: It is not very often that an action takes the system to the best possible state with $50\%$ probability and to the worst possible state with $50\%$ probability.

Most PAC exploration algorithms take an average over samples. By contrast, the algorithm presented in this paper splits samples into sets, takes the average over each set, and returns the median of the averages. This seemingly simple trick (known as the median trick [1]), allows us to derive sample complexity bounds that depend on the variance of the Bellman operator rather than $Q_{\max}^2$. Addi-

tionally, our algorithm (Median-PAC) is the first reinforcement learning algorithm with theoretical guarantees that allows for unbounded rewards.[1]

Not only does Median-PAC offer significant sample complexity savings in the case when the variance of the Bellman operator is low, but even in the worst case (the variance of the Bellman operator is bounded above by $\frac{Q_{\max}^2}{4}$) our bounds match the best, published PAC bounds. Note that Median-PAC does not require the variance of the Bellman operator to be known in advance. Our bounds show that there is an inverse relationship between the (possibly unknown) variance of the Bellman operator and Median-PAC's performance. This is to the best of our knowledge not only the first application of the median of means in PAC exploration, but also the first application of the median of means in reinforcement learning in general.

Contrary to recent work which has exploited variance in Markov decision processes to improve PAC bounds [7, 3], Median-PAC makes no assumptions about the number of possible next-states from every state-action (it does not even require the number of possible next states to be finite), and as a result it is easily extensible to the continuous state, concurrent MDP, and delayed update settings [12].

## 2 Background, notation, and definitions

In the following, important symbols and terms will appear in **bold** when first introduced. Let $\mathcal{X}$ be the domain of $x$. Throughout this paper, $\forall x$ will serve as a shorthand for $\forall x \in \mathcal{X}$. In the following $\boldsymbol{s}, \boldsymbol{\bar{s}}, \boldsymbol{\tilde{s}}, \boldsymbol{s'}$ are used to denote various states, and $\boldsymbol{a}, \boldsymbol{\bar{a}}, \boldsymbol{\tilde{a}}, \boldsymbol{a'}$ are used to denote actions.

A *Markov Decision Process* (MDP) [13] is a 5-tuple $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$, where $\mathcal{S}$ is the state space of the process, $\boldsymbol{A}$ is the action space[2], $\boldsymbol{P}$ is a Markovian transition model $\big(p(s'|s, a)$ denotes the probability of a transition to state $s'$ when taking action $a$ in state $s\big)$, $\boldsymbol{R}$ is a reward function $\big(R(s, a, s')$ is the reward for taking action $a$ in state $s$ and transitioning to state $s'\big)$, and $\boldsymbol{\gamma} \in [0, 1)$ is a discount factor for future rewards. A *deterministic policy* $\boldsymbol{\pi}$ is a mapping $\pi : \mathcal{S} \mapsto \mathcal{A}$ from states to actions; $\pi(s)$ denotes the action choice in state $s$. The value $\boldsymbol{V^{\pi}(s)}$ of state $s$ under policy $\pi$ is defined as the expected, accumulated, discounted reward when the process begins in state $s$ and all decisions are made according to policy $\pi$. There exists an optimal policy $\boldsymbol{\pi^*}$ for choosing actions which yields the optimal value function $V^*(s)$, defined recursively via the Bellman optimality equation $\boldsymbol{V^*(s)} = \max_a \{\sum_{s'} p(s'|s, a) (R(s, a, s') + \gamma V^*(s'))\}$. Similarly, the value $\boldsymbol{Q^{\pi}(s, a)}$ of a state-action $(s, a)$ under policy $\pi$ is defined as the expected, accumulated, discounted reward when the process begins in state $s$ by taking action $a$ and all decisions thereafter are made according to policy $\pi$. The Bellman optimality equation for $Q$ becomes $\boldsymbol{Q^*(s, a)} = \sum_{s'} p(s'|s, a) (R(s, a, s') + \gamma \max_{a'} \{Q^*(s', a')\})$. For a fixed policy $\pi$ the Bellman operator for $Q$ is defined as $\boldsymbol{B^{\pi}Q(s, a)} = \sum_{s'} p(s'|s, a)\big(R(s, a, s') + \gamma Q(s', \pi(s'))\big)$. In reinforcement learning (RL) [15], a learner interacts with a stochastic process modeled as an MDP and typically observes the state and immediate reward at every step; however, the transition model $P$ and reward function $R$ are not known. The goal is to learn a near optimal policy using experience collected through interaction with the process. At each step of interaction, the learner observes the current state $s$, chooses an action $a$, and observes the reward received $r$, and resulting next state $s'$, essentially sampling the transition model and reward function of the process. Thus experience comes in the form of $(s, a, r, s')$ samples.

We assume that all value functions $Q$ live in a complete metric space.

**Definition 2.1.** $\boldsymbol{Q_{\mathbf{max}}}$ *denotes an upper bound on the expected, accumulated, discounted reward from any state-action under any policy.*

We require that $Q_{\min}$, the minimum expected, accumulated, discounted reward from any state-action under any policy is bounded, and in order to simplify notation we also assume without loss of

---

[1]Even though domains with truly unbounded rewards are not common, many domains exist for which infrequent events with extremely high (winning the lottery) or extremely low (nuclear power-plant meltdown) rewards exist. Algorithms whose sample complexity scales with the highest magnitude event are not well suited to such domains.

[2]For simplicity of exposition we assume that the same set of actions is available at every state. Our results readily extend to the case where the action set can differ from state to state.

generality that it is bounded below by 0. If $Q_{\min} < 0$, this assumption is easy to satisfy in all MDPs for which $Q_{\min}$ is bounded by simply shifting the reward space by $(\gamma - 1)Q_{\min}$.

There have been many definitions of **sample complexity** in RL. In this paper we will be using the following [12]:

**Definition 2.2.** *Let $(s_1, s_2, s_3, \dots)$ be the random path generated on some execution of $\pi$, where $\pi$ is an arbitrarily complex, possibly non-stationary, possibly history dependent policy (such as the policy followed by an exploration algorithm). Let $\epsilon$ be a positive constant, $T$ the (possibly infinite) set of time steps for which $V^{\pi}(s_t) < V^*(s_t) - \epsilon$, and define[3]*

$$\epsilon_e(t) = V^*(s_t) - V^{\pi}(s_t) - \epsilon, \ \forall \, t \in T.$$
$$\epsilon_e(t) = 0, \ \forall \, t \notin T.$$

*The Total Cost of Exploration (**TCE**) is defined as the undiscounted infinite sum $\sum_{t=0}^{\infty} \epsilon_e(t)$.*

"Number of suboptimal steps" bounds follow as a simple corollary of TCE bounds.

We will be using the following definition of efficient PAC exploration [14]:

**Definition 2.3.** *An algorithm is said to be **efficient PAC-MDP** (Probably Approximately Correct in Markov Decision Processes) if, for any $\epsilon > 0$ and $0 < \delta < 1$, its sample complexity, its per-timestep computational complexity, and its space complexity, are less than some polynomial in the relevant quantities $(S, A, \frac{1}{\epsilon}, \frac{1}{\delta}, \frac{1}{1-\gamma})$, with probability at least $1 - \delta$.*

## 3 The median of means

Before we present Median-PAC we will demonstrate the usefulness of the median of means with a simple example. Suppose we are given $n$ independent samples from a random variable $X$ and we want to estimate its mean. The types of guarantees that we can provide about how close that estimate will be to the expectation, will depend on what knowledge we have about the variable, and on the method we use to compute the estimate. The main question of interest in our work is how many samples are needed until our estimate is $\epsilon$-close to the expectation with probability at least $1 - \delta$.

Let the expectation of $X$ be $E[X] = \mu$ and its variance $var[X] = \sigma^2$. Cantelli's inequality tells us that: $P(X - \mu \geq \epsilon) \leq \frac{\sigma^2}{\sigma^2 + \epsilon^2}$ and $P(X - \mu \leq -\epsilon) \leq \frac{\sigma^2}{\sigma^2 + \epsilon^2}$. Let $X_i$ be a random variable describing the value of the $i$-th sample, and define $X' = \frac{X_1 + X_2 + \dots + X_n}{n}$. We have that $E[X'] = \mu$ and $var[X'] = \frac{\sigma^2}{n}$. From Cantelli's inequality we have that $P(X' - \mu \geq \epsilon) \leq \frac{\sigma^2}{\sigma^2 + n\epsilon^2}$ and $P(X' - \mu \leq -\epsilon) \leq \frac{\sigma^2}{\sigma^2 + n\epsilon^2}$. Solving for $n$ we have that we need at most $n = \frac{(1-\delta)\sigma^2}{\delta\epsilon^2} = O\left(\frac{\sigma^2}{\delta\epsilon^2}\right)$ samples until our estimate is $\epsilon$-close to the expectation with probability at least $1 - \delta$. In RL, it is common to apply a union bound over the entire state-action space in order to prove uniformly good approximation. This means that $\delta$ has to be small enough that even when multiplied with the number of state-actions, it yields an acceptably low probability of failure. The most significant drawback of the bound above is that it grows very quickly as $\delta$ becomes smaller. Without further assumptions one can show that the bound above is tight for the average estimator.

If we know that $X$ can only take values in a bounded range $a \leq X \leq b$, Hoeffding's inequality tells us that $P(X' - \mu \geq \epsilon) \leq e^{-\frac{2n\epsilon^2}{(b-a)^2}}$ and $P(X' - \mu \leq -\epsilon) \leq e^{-\frac{2n\epsilon^2}{(b-a)^2}}$. Solving for $n$ we have that $n = \frac{(b-a)^2 \ln \frac{1}{\delta}}{2\epsilon^2}$ samples suffice to guarantee that our estimate is $\epsilon$-close to the expectation with probability at least $1 - \delta$. Hoeffding's inequality yields a much better bound with respect to $\delta$, but introduces a quadratic dependence on the range of values that the variable can take. For long planning horizons (discount factor close to 1) and/or large reward magnitudes, the range of possible $Q$-values can be very large, much larger than the variance of individual state-actions.

We can get the best of both worlds by using a more sophisticated estimator. Instead of taking the average over $n$ samples, we will split them into $k_m = \frac{n\epsilon^2}{4\sigma^2}$ sets of $\frac{4\sigma^2}{\epsilon^2}$ samples each,[4] compute the

---

[3]Note that $V^{\pi}(s_t)$ denotes the expected, discounted, accumulated reward of the arbitrarily complex policy $\pi$ from state $s_t$ at time $t$, rather than the expectation of some stationary snapshot of $\pi$.

[4]The number of samples per set was chosen so as to minimize the constants in the final bound.

average over each set, and then take the median of the averages. From Cantelli's inequality we have that with probability at least $\frac{4}{5}$, each one of the sets will not underestimate, or overestimate the mean $\mu$ by more than $\epsilon$. Let $f^-$ be the function that counts the number of sets that underestimate the mean by more than $\epsilon$, and $f^+$ the function that counts the number of sets that overestimate the mean by more than $\epsilon$. From McDiarmid's inequality [9] we have that $P\left(f^- \geq \frac{k_m}{2}\right) \leq e^{-\frac{2\left(\frac{3k_m}{10}\right)^2}{k_m}}$ and $P\left(f^+ \geq \frac{k_m}{2}\right) \leq e^{-\frac{2\left(\frac{3k_m}{10}\right)^2}{k_m}}$. Solving for $n$ we have that $n = \frac{\frac{200}{9}\sigma^2 \ln\left(\frac{1}{\delta}\right)}{\epsilon^2} \approx \frac{22.22\sigma^2 \ln\left(\frac{1}{\delta}\right)}{\epsilon^2}$ samples suffice to guarantee that our estimate is $\epsilon$-close to the expectation with probability at least $1 - \delta$. The median of means offers logarithmic dependence on $\frac{1}{\delta}$, independence from the range of values that the variables in question can take (even allowing for them to be infinite), and can be computed efficiently. The median of means estimator only requires a finite variance and the existence of a mean. No assumptions (including boundedness) are made on higher moments.

## 4 Median PAC exploration

---
**Algorithm 1** Median-PAC

---
1: Inputs: start state $s$, discount factor $\gamma$, max number of samples $k$, number of sets $k_m$, and acceptable error $\epsilon_a$.
2: Initialize sample sets $u_{new}(s,a) = \emptyset, u(s,a) = \emptyset \ \forall \ (s,a)$. ($|u(s,a)|$ denotes the number of samples in $u(s,a)$)
3: Set $\epsilon_b = \epsilon_a \sqrt{k}$, and initialize value function $\tilde{Q}(s,a) = Q_{\max} \ \forall \ (s,a)$.
4: **loop**
5:     Perform action $a = \arg\max_{\tilde{a}} \tilde{Q}(s, \tilde{a})$
6:     Receive reward $r$, and transition to state $s'$.
7:     **if** $|u(s,a)| < k$ **then**
8:         Add $(s, a, r, s')$ to $u_{new}(s,a)$.
9:         **if** $|u_{new}(s,a)| > |u(s,a)|$ and $|u_{new}(s,a)| = 2^i k_m$, where $i \geq 0$ is an integer **then**
10:           $u(s,a) = u_{new}(s,a)$
11:           $u_{new}(s,a) = \emptyset$
12:         **end if**
13:         **while** $max_{(s,a)}(\tilde{B}\tilde{Q}(s,a) - \tilde{Q}(s,a)) > \epsilon_a$ **or** $max_{(s,a)}(\tilde{Q}(s,a) - \tilde{B}\tilde{Q}(s,a)) > \epsilon_a$ **do**
14:           Set $\tilde{Q}(s,a) = \tilde{B}\tilde{Q}(s,a) \ \forall \ (s,a)$.
15:         **end while**
16:     **end if**
17: **end loop**
18: **function** $\tilde{B}\tilde{Q}(s,a)$
19:     **if** $|u(s,a)| \geq k_m$ **then**
20:         Let $(s, a, r_i, s'_i)$ be the $i$-th sample in $u(s,a)$.
21:         **for** $j = 1$ **to** $k_m$ **do**
22: $$g(j) = \sum_{i=1+(j-1)\frac{|u(s,a)|}{k_m}}^{j\frac{|u(s,a)|}{k_m}} \left( r_i + \gamma \max_{\bar{a}} \tilde{Q}(s'_i, \bar{a}) \right)$$
23:         **end for**
24:         **return** $\min\left\{ Q_{\max}, \frac{\epsilon_b}{\sqrt{|u(s,a)|}} + \frac{k_m \, median\{g(1),...g(k_m)\}}{|u(s,a)|} \right\}$
25:     **else**
26:         **return** $Q_{\max}$.
27:     **end if**
28: **end function**

---

Algorithm 1 has three parameters that can be set by the user:

- $k$ is the maximum number of samples per state-action. As we will show, higher values for $k$ lead to increased sample complexity but better approximation.

- $\epsilon_a$ is an "acceptable error" term. Since Median-PAC is based on value iteration (lines 13 through 15) we specify a threshold after which value iteration should terminate. Value

iteration is suspended when the max-norm of the difference between Bellman backups is no larger than $\epsilon_a$.

- Due to the stochasticity of Markov decision processes, Median-PAC is only guaranteed to achieve a particular approximation quality with some probability. $k_m$ offers a trade-off between approximation quality and the probability that this approximation quality is achieved. For a fixed $k$ smaller values of $k_m$ offer potentially improved approximation quality, while larger values offer a higher probability of success. For simplicity of exposition our analysis requires that $k = 2^i k_m$ for some integer $i$. If $k_m \geq \left\lceil \frac{50}{9} \ln \frac{4 \log_2 \frac{4Q_{\max}^2}{\epsilon_a^2} |SA|^2}{\delta} \right\rceil$ the probability of failure is bounded above by $\delta$.

Like most modern PAC exploration algorithms, Median-PAC is based on the principle of optimism in the face of uncertainty. At every step, the algorithm selects an action greedily based on the current estimate of the $Q$-value function $\tilde{Q}$. The value function is optimistically initialized to $Q_{\max}$, the highest value that any state-action can take. If $k$ is set appropriately (see theorem 5.4), the value function is guaranteed to remain approximately optimistic (approximately represent the most optimistic world consistent with the algorithm's observations) with high probability.

We would like to draw the reader's attention to two aspects of Median-PAC, both in the way Bellman backups are computed: 1) Instead of taking a simple average over sample values, Median-PAC divides them into $k_m$ sets, computes the mean over each set, and takes the median of means. 2) Instead of using all the samples available for every state-action, Median-PAC uses samples in batches of a power of 2 times $k_m$ (line 9). The reasoning behind the first choice follows from the discussion above: using the median of means will allow us to show that Median-PAC's complexity scales with the variance of the Bellman operator (see definition 5.1) rather than $Q_{\max}^2$. The reasoning behind using samples in batches of increasing powers of 2 is more subtle. A key requirement in the analysis of our algorithm is that samples belonging to the same state-action are independent. While the outcome of sample $i$ does not provide information about the outcome of sample $j$ if $i < j$ (from the Markov property), the fact that $j$ samples exist can reveal information about the outcome of $i$. If the first $i$ samples led to a severe underestimation of the value of the state-action in question, it is likely that $j$ samples would never have been collected. The fact that they did gives us some information about the outcome of the first $i$ samples. Using samples in batches, and discarding the old batch when a new batch becomes available, ensures that the outcomes of samples within each batch are independent from one another.

## 5 Analysis

**Definition 5.1.** $\sigma$ *is the minimal constant satisfying*

$$\forall (s, a, \pi^{\tilde{Q}}, \tilde{Q}), \sqrt{\sum_{s'} p(s'|s, a)\left( R(s, a, s') + \gamma \tilde{Q}(s', \pi^{\tilde{Q}}(s')) - B^{\pi^{\tilde{Q}}} \tilde{Q}(s, a) \right)^2} \leq \sigma,$$

*where $\forall \tilde{Q}$ refers to any value function produced by Median-PAC, rather than any conceivable value function (similarly $\pi^{\tilde{Q}}$ refers to any greedy policy over $\tilde{Q}$ followed during the execution of Median-PAC rather than any conceivable policy).*

In the following we will call $\sigma^2$ the variance of the Bellman operator. Note that the variance of the Bellman operator is not the same as the variance, or stochasticity in the transition model of an MDP. A state-action can be highly stochastic (lead to many possible next states), yet if all the states it transitions to have similar values, the variance of its Bellman operator will be small.

From Lemmas 5.2, 5.3, and theorem 5.4 below, we have that Median-PAC is efficient PAC-MDP.

**Lemma 5.2.** *The space complexity of algorithm 1 is $O\left(k|S||A|\right)$.*

*Proof.* Follows directly from the fact that at most $k$ samples are stored per state-action. $\square$

**Lemma 5.3.** *The per step computational complexity of algorithm 1 is bounded above by*

$$O\left( \frac{k|S||A|^2}{1 - \gamma} \ln \frac{Q_{\max}}{\epsilon_a} \right).$$

*Proof.* The proof of this lemma is deferred to the appendix. □

Theorem 5.4 below is the main theorem of this paper. It decomposes errors into the following three sources:

1. $\epsilon_a$ is the error caused by the fact that we are only finding an $\epsilon_a$-approximation, rather than the true fixed point of the approximate Bellman operator $\tilde{B}$, and the fact that we are using only a finite set of samples (at most $k$) to compute the median of the means, thus we only have an estimate.

2. $\epsilon_u$ is the error caused by underestimating the variance of the MDP. When $k$ too small and Median-PAC fails to be optimistic, $\epsilon_u$ will be non-zero. $\epsilon_u$ is a measure of how far Median-PAC is from being optimistic (follow the greedy policy over the value function of the most optimistic world consistent with its observations).

3. Finally, $\epsilon_e(t)$ is the error caused by the fact that at time $t$ there may exist state-actions that do not yet have $k$ samples.

**Theorem 5.4.** *Let* $(s_1, s_2, s_3, \dots)$ *be the random path generated on some execution of Median-PAC, and* $\tilde{\pi}$ *be the (non-stationary) policy followed by Median-PAC. Let* $\epsilon_u = \max\{0, \sigma\sqrt{4k_m} - \epsilon_a\sqrt{k}\}$, *and* $\epsilon_a$ *be defined as in algorithm 1. If* $k_m = \left\lceil \frac{50}{9} \ln \frac{4 \log_2 \frac{4Q_{\max}^2 |SA|^2}{\epsilon_a^2}}{\delta} \right\rceil$,

$\frac{2\left\lceil \frac{1}{1-\gamma} \ln \frac{(1-\gamma)Q_{\max}}{\epsilon_a} \right\rceil 2 \ln \frac{\log_2 \frac{2k}{k_m}}{\delta}}{k_m |SA| + 1} < 1$, *and* $k = 2^i k_m$ *for some integer* $i$, *then with probability at least* $1 - \delta$, *for all* $t$

$$V^*(s_t) - V^{\tilde{\pi}}(s_t) \le \frac{2\epsilon_u + 5\epsilon_a}{1 - \gamma} + \epsilon_e(t), \tag{1}$$

*where*

$$\sum_{t=0}^{\infty} \epsilon_e(t) < c_0 \left( \left( 2k_m + \log_2 \frac{2k}{k_m} \right) Q_{\max} + \epsilon_a k \left( 8 + \frac{8}{\sqrt{2}} \right) \right), \tag{2}$$

*and*

$$c_0 = \frac{(|SA| + 1)\left(1 + \log_2 \left\lceil \frac{1}{1-\gamma} \ln \frac{(1-\gamma)Q_{\max}}{\epsilon_a} \right\rceil\right) \left\lceil \frac{1}{1-\gamma} \ln \frac{(1-\gamma)Q_{\max}}{\epsilon_a} \right\rceil}{1 - \sqrt{\frac{2\left\lceil \frac{1}{1-\gamma} \ln \frac{(1-\gamma)Q_{\max}}{\epsilon_a} \right\rceil 2 \ln \frac{\log_2 \frac{2k}{k_m}}{\delta}}{k_m |SA| + 1}}}.$$

*If* $k = 2^i k_m$ *where* $i$ *is the smallest integer such that* $2^i \ge \frac{4\sigma^2}{\epsilon_a^2}$, *and* $\epsilon_0 = (1 - \gamma)\epsilon_a$, *then with probability at least* $1 - \delta$, *for all* $t$

$$V^*(s_t) - V^{\tilde{\pi}}(s_t) \le \epsilon_0 + \epsilon_e(t), \tag{3}$$

*where*[5]

$$\sum_{t=0}^{\infty} \epsilon_e(t) \approx \tilde{O}\left( \left( \frac{\sigma^2}{\epsilon_0 (1-\gamma)^2} + \frac{Q_{\max}}{1 - \gamma} \right) |SA| \right). \tag{4}$$

*Note that the probability of success holds for all timesteps simultaneously, and* $\sum_{t=0}^{\infty} \epsilon_e(t)$ *is an undiscounted infinite sum.*

*Proof.* The detailed proof of this theorem is deferred to the appendix. Here we provide a proof sketch:

The non-stationary policy of the algorithm can be broken up into fixed policy (and fixed approximate value function) segments. The first step in proving theorem 5.4 is to show that the Bellman error of each state-action at a particular fixed approximate value function segment is acceptable with respect to the number of samples currently available for that state-action with high probability. We use Cantelli's and McDiarmid's inequalities to prove this point. This is where the median of means

---

[5] $f(n) = \tilde{O}(g(n))$ is a shorthand for $f(n) = O(g(n) \log^c g(n))$ for some constant $c$.

becomes useful, and the main difference between our work and earlier work. We then combine the result from the median of means, the fact that there are only a small number of possible policy and approximate value function changes that can happen during the lifetime of the algorithm, and the union bound, to prove that the Bellman error of all state-actions during all timesteps is acceptable with high probability. We subsequently prove that due to the optimistic nature of Median-PAC, at every time-step it will either perform well, or learn something new about the environment with high probability. Since there is only a finite number of things it can learn, the total cost of exploration for Median-PAC will be small with high probability. □

A typical "number of suboptimal steps" sample complexity bound follows as a simple corollary of theorem 5.4. If the total cost of exploration is $\sum_{t=0}^{\infty} \epsilon_e(t)$ for an $\epsilon_0$-optimal policy, there can be no more than $\frac{\sum_{t=0}^{\infty} \epsilon_e(t)}{\epsilon_1}$ steps that are more than $(\epsilon_0 + \epsilon_1)$-suboptimal.

Note that the sample complexity of Median-PAC depends log-linearly on $Q_{\max}$, which can be finite even if $R_{\max}$ is infinite. Consider for example an MDP for which the reward at every state-action follows a Gaussian distribution (for discrete MDPs this example requires rewards to be stochastic, while for continuous MDPs rewards can be a deterministic function of state-action-nextstate since there can be an infinite number of possible nextstates for every state-action). If the mean of the reward for every state-action is bounded above by $c$, $Q_{\max}$ is bounded above by $\frac{c}{1-\gamma}$, even though $R_{\max}$ is infinite.

As we can see from theorem 5.4, apart from being the first PAC exploration algorithm that can be applied to MDPs with unbounded rewards, Median-PAC offers significant advantages over the current state of the art for MDPs with bounded rewards. Until recently, the algorithm with the best known sample complexity for the discrete state-action setting was MORMAX, an algorithm by Szita and Szepesvári [16]. Theorem 5.4 offers an improvement of $\frac{1}{(1-\gamma)^2}$ even in the worst case, and trades a factor of $Q_{\max}^2$ for a (potentially much smaller) factor of $\sigma^2$. A recent algorithm by Pazis and Parr [12] currently offers the best known bounds for PAC exploration without additional assumptions on the number of states that each action can transition to. Compared to that work we trade a factor of $Q_{\max}^2$ for a factor of $\sigma^2$.

## 5.1 Using Median-PAC when $\sigma$ is not known

In many practical situations $\sigma$ will not be known. Instead the user will have a fixed exploration cost budget, a desired maximum probability of failure $\delta$, and a desired maximum error $\epsilon_a$. Given $\delta$ we can solve for the number of sets as $k_m = \left\lceil \frac{50}{9} \ln \frac{4 \log_2 \frac{4Q_{\max}^2 |SA|^2}{\epsilon_a^2}}{\delta} \right\rceil$, at which point all variables in equation 2 except for $k$ are known, and we can solve for $k$. When the sampling budget is large enough such that $k \geq \frac{4\sigma^2 k_m}{\epsilon_a^2}$, then $\epsilon_u$ in equation 1 will be zero. Otherwise $\epsilon_u = \sigma\sqrt{4k_m} - \epsilon_a\sqrt{k}$.

## 5.2 Beyond the discrete state-action setting

Recent work has extended PAC exploration to the continuous state [11] concurrent exploration [4] and delayed update [12] settings. The goal in the concurrent exploration setting is to explore in multiple identical or similar MDPs and incur low aggregate exploration cost over all MDPs. For a concurrent algorithm to offer an improvement over non-concurrent exploration, the aggregate cost must be lower than the cost of non-concurrent exploration times the number of tasks. The delayed update setting takes into account the fact that in real world domains, reaching a fixed point after collecting a new sample can take longer that the time between actions. Contrary to other work that has exploited the variance of MDPs to improve bounds on PAC exploration [7, 3] our analysis does not make assumptions about the number of possible next states from a given action. As such, Median-PAC and its bounds are easily extensible to the continuous state, concurrent exploration, delayed update setting. Replacing the average over samples in an approximation unit with the median of means over samples in an approximation unit in the algorithm of Pazis and Parr [12], improves their bounds (which are the best published bounds for PAC exploration in these settings) by $(R_{\max} + \gamma Q_{\max})^2$ while introducing a factor of $\sigma^2$.

# 6 Experimental evaluation

We compared Median-PAC against the algorithm of Pazis and Parr [12] on a simple 5 by 5 gridworld (see appendix for more details). The agent has four actions: move one square up, down, left, or right. All actions have a $1\%$ probability of self-transition with a reward of 100. Otherwise the agent moves in the chosen direction and receives a reward of 0, unless its action causes it to land on the top-right corner, in which case it receives a reward of 1. The world wraps around and the agent always starts at the center. The optimal policy for this domain is to take the shortest path to the top-right corner if at a state other than the top-right corner, and take any action while at the top-right corner.

While the probability of any individual sample being a self-transition is small, unless the number of samples per state-action is very large, the probability that there will exist at least one state-action with significantly more than $\frac{1}{100}$ sampled self-transitions is high. As a result, the naive average algorithm frequently produced a policy that maximized the probability of encountering state-actions with more than $\frac{1}{100}$ sampled self-transitions. By contrast, it is far less likely that there will exist a state-action for which at least half of the sets used by the median of means have more than $\frac{1}{100}$ sampled self-transitions. Median-PAC was able to consistently find the optimal policy.

# 7 Related Work

Maillard, Mann, and Mannor [8] present the distribution norm, a measure of hardness of an MDP. Similarly to our definition of the variance of the Bellman operator, the distribution norm does not directly depend on the stochasticity of the underlying transition model. It would be interesting to see if the distribution norm (or a similar concept) can be used to improve PAC exploration bounds for "easy" MDPs.

While to the best our knowledge our work is the first in PAC exploration for MDPs that introduces a measure of hardness for MDPs (the variance of the Bellman operator), measures of hardness have been previously used in regret analysis [6]. Such measures include the diameter of an MDP [6], the one way diameter [2], as well as the span [2]. These measures express how hard it is to reach any state of an MDP from any other state. A major advantage of sample complexity over regret is that finite diameter is not required to prove PAC bounds. Nevertheless, if introducing a requirement for a finite diameter could offer drastically improved PAC bounds, it may be worth the trade-off for certain classes of problems. Note that variance and diameter of an MDP appear to be orthogonal. One can construct examples of arbitrary diameter and then manipulate the variance by changing the reward function and/or discount factor.

Another measure of hardness which was recently introduced in regret analysis is the Eluder dimension. Osband and Van Roy [10] show that if an MDP can be parameterized within some known function class, regret bounds that scale with the dimensionality, rather than cardinality of the underlying MDP can be obtained. Like the diameter, the Eluder dimension appears to be orthogonal to the variance of the Bellman operator, potentially allowing for the two concepts to be combined.

Lattimore and Hutter [7] have presented an algorithm that can match the best known lower bounds for PAC exploration up to logarithmic factors for the case of discrete MDPs where every state-action can transition to at most two next states.

To the best of our knowledge there has been no work in learning with unbounded rewards. Harrison [5] has examined the feasibility of planning with unbounded rewards.

## Acknowledgments

# References

[1] N Alon, Y Matias, and M Szegedy. The space complexity of approximating the frequency moments. *Journal of Computer and System Sciences - JCSS (special issue of selected papers from STOC'96)*, 58:137–147, 1999.

[2] Peter L. Bartlett and Ambuj Tewari. REGAL: A regularization based algorithm for reinforcement learning in weakly communicating MDPs. In *Proceedings of the 25th Conference on Uncertainty in Artificial Intelligence (UAI2009)*, pages 35–42, June 2009.

[3] Christoph Dann and Emma Brunskill. Sample complexity of episodic fixed-horizon reinforcement learning. *Advances in Neural Information Processing Systems*, 2015.

[4] Zhaohan Guo and Emma Brunskill. Concurrent PAC RL. In *AAAI Conference on Artificial Intelligence*, pages 2624–2630, 2015.

[5] J. Michael Harrison. Discrete dynamic programming with unbounded rewards. *The Annals of Mathematical Statistics*, 43(2):636–644, 04 1972.

[6] Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11:1563–1600, August 2010.

[7] Tor Lattimore and Marcus Hutter. PAC bounds for discounted MDPs. In *Proceedings of the 23th International Conference on Algorithmic Learning Theory*, volume 7568 of *Lecture Notes in Computer Science*, pages 320–334. Springer Berlin / Heidelberg, 2012.

[8] Odalric-Ambrym Maillard, Timothy A Mann, and Shie Mannor. How hard is my MDP?" the distribution-norm to the rescue". In *Advances in Neural Information Processing Systems 27*, page 1835–1843. 2014.

[9] C. McDiarmid. On the method of bounded differences. In *Surveys in Combinatorics*, number 141 in London Mathematical Society Lecture Note Series, pages 148–188. Cambridge University Press, August 1989.

[10] Ian Osband and Benjamin Van Roy. Model-based reinforcement learning and the eluder dimension. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 1466–1474. 2014.

[11] Jason Pazis and Ronald Parr. PAC optimal exploration in continuous space Markov decision processes. In *AAAI Conference on Artificial Intelligence*, pages 774–781, July 2013.

[12] Jason Pazis and Ronald Parr. Efficient PAC-optimal exploration in concurrent, continuous state MDPs with delayed updates. In *AAAI Conference on Artificial Intelligence*, February 2016.

[13] Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wiley-Interscience, April 1994.

[14] Alexander L. Strehl, Lihong Li, and Michael L. Littman. Reinforcement learning in finite MDPs: PAC analysis. *Journal of Machine Learning Research*, 10:2413–2444, December 2009.

[15] Richard Sutton and Andrew Barto. *Reinforcement Learning: An Introduction*. The MIT Press, Cambridge, Massachusetts, 1998.

[16] Istvan Szita and Csaba Szepesvári. Model-based reinforcement learning with nearly tight exploration complexity bounds. In *International Conference on Machine Learning*, pages 1031–1038, 2010.

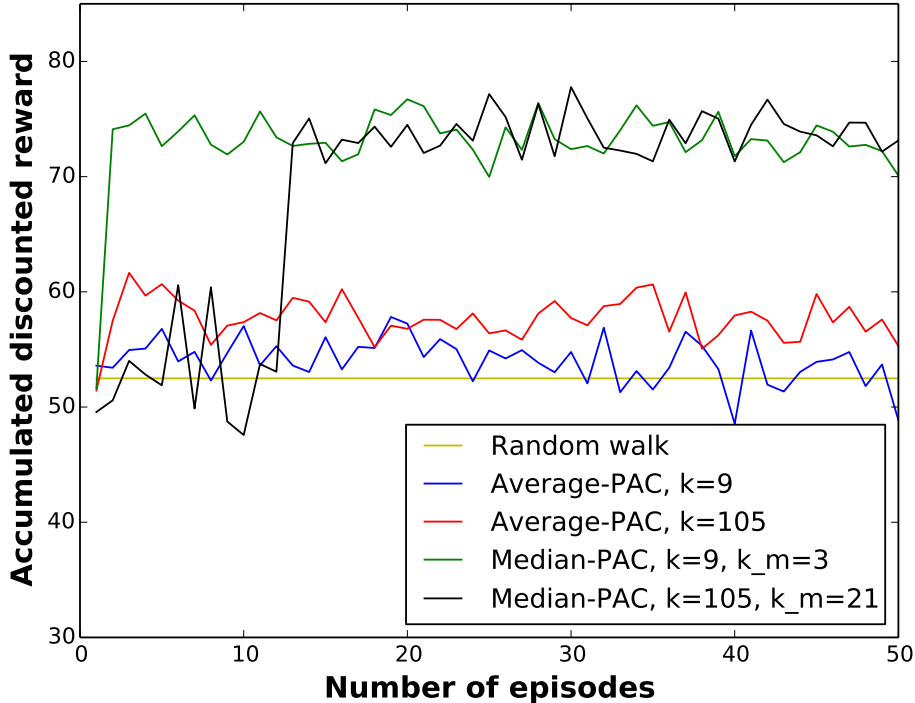# 8 Appendix A: Experimental evaluation



Figure 1: Accumulated discounted reward as a function of the number of episodes for a random walk, the algorithm of Pazis and Parr [12], and Median-PAC on a simple gridworld. Each plot represents an average over 1000 independent repetitions.

The discount factor for the gridworld described in section 6 was set to $0.98$, and every episode was 1000 steps long. We used modified versions of both learning algorithms that accumulate samples rather than using them in batches and discarding the old, smaller batch once a new batch has been collected. The algorithm of Pazis and Parr [12] (Average-PAC), was allowed allowed 1000 iterations of value iteration after each sample was added. Median-PAC was allowed 1000 iterations of value iteration every time the $i * k_m$-th sample was added to a state action, where $i > 0$ is an integer. $\epsilon_b$ was set to $0.01Q_{max}$ for both algorithms (since both algorithms truncate state-action values to $Q_{max}$, setting $\epsilon_b$ close to $Q_{max}$ for small values of $k$ saturates the value function). Notice that Median-PAC for $k = 105$ and $k_m = 21$ takes longer to achieve good performance than for $k = 9$ and $k_m = 3$. This is because for $k_m = 21$ the value of every state-action state is $Q_{max}$ until at least 21 samples have been collected.

# 9 Appendix B: Analysis

Before we prove lemma 5.3 and theorem 5.4 we have to introduce a few supporting definitions and lemmas.

**Definition 9.1.** *Let* $|u(s,a)| = 2^i k_m$ *for some* $i \in \{1, 2, \dots\}$. *The function* $\boldsymbol{F^\pi(Q, u(s, a))}$ *is defined as*

$$F^\pi(Q, u(s, a)) = \frac{\epsilon_b}{\sqrt{|u(s, a)|}} + median\{G^\pi(Q, u(s, a), 1),$$

$$\dots,$$

$$G^\pi(Q, u(s, a), k_m)\},$$

*where*

$$G^\pi(Q, u(s, a), j) = \frac{k_m}{|u(s, a)|} \sum_{i=1+(j-1)\frac{|u(s,a)|}{k_m}}^{j\frac{|u(s,a)|}{k_m}} \left( r_i + \gamma Q(s_i', \pi(s_i')) \right),$$

*and* $(s, a, r_i, s_i')$ *is the* $i$*-th sample in* $u(s, a)$. *We will use* $\boldsymbol{F(Q, u(s, a))}$ *to denote* $F^{\pi^Q}(Q, u(s, a))$.

$F^\pi$ splits the samples in $u(s, a)$ into $k_m$ groups, computes the average of the sample values in each group, and returns the median of the averages.

**Definition 9.2.** *For state-action* $(s, a)$, *the approximate optimistic Bellman operator* $\boldsymbol{\tilde{B}^\pi}$ *for policy* $\pi$ *is defined as*

$$\tilde{B}^\pi Q(s, a) = \min\{Q_{\max}, F^\pi(Q, u(s, a))\}.$$

*We will use* $\boldsymbol{\tilde{B}Q(s, a)}$ *to denote* $\tilde{B}^{\pi^Q} Q(s, a)$. *When* $|u(s, a)| = 0$, $\tilde{B}^\pi Q(s, a) = Q_{\max}$.

The approximate optimistic Bellman operator is applied to the approximate value function on line 14 of the algorithm.

**Lemma 9.3.** $\tilde{B}$ *is a* $\gamma$-*contraction in maximum norm.*

*Proof.* Suppose $||Q_1 - Q_2||_\infty = \epsilon$. For any $(s, a)$ we have

$$\begin{aligned}
\tilde{B}Q_1(s, a) &= \min\{Q_{\max}, F(Q_1, u(s, a))\} \\
&\leq \min\{Q_{\max}, F(Q_2, u(s, a)) + \gamma\epsilon\} \\
&\leq \gamma\epsilon + \min\{Q_{\max}, F(Q_2, u(s, a))\} \\
&= \gamma\epsilon + \tilde{B}Q_2(s, a) \\
&\Rightarrow \tilde{B}Q_1(s, a) \leq \gamma\epsilon + \tilde{B}Q_2(s, a).
\end{aligned}$$

Similarly we have that $\tilde{B}Q_2(s, a) \leq \gamma\epsilon + \tilde{B}Q_1(s, a)$ which completes our proof. $\qquad\square$

**Lemma 9.4.** *Let* $\sigma$ *be defined as in Definition* 5.1. *For a fixed* $\tilde{Q}$ *and fixed* $(s, a)$ *such that* $|u(s, a)| > 0$

$$P\left(G^\pi(\tilde{Q}, u(s, a), j) - B^\pi \tilde{Q}(s, a) \leq -\frac{\sigma\sqrt{4k_m}}{\sqrt{|u(s, a)|}}\right) \leq \frac{1}{5},$$

*and*

$$P\left(G^\pi(\tilde{Q}, u(s, a), j) - B^\pi \tilde{Q}(s, a) \geq \frac{\sigma\sqrt{4k_m}}{\sqrt{|u(s, a)|}}\right) \leq \frac{1}{5}.$$

*Proof.* From Definition 9.1 we have that

$$B^\pi \tilde{Q}(s, a) = \mathrm{E}\left[G^\pi(\tilde{Q}, u(s, a), j)\right],$$

where the expectation is over the next-states that samples in $u(s, a)$ used by $G^\pi$ land on.

Let $Y$ be the set of $\frac{|u(s,a)|}{k_m}$ samples used by $G^\pi(\tilde{Q}, u(s, a), j)$ at $(s, a)$. Define $Z_1, \dots Z_{\frac{|u(s,a)|}{k_m}}$ to be random variables, one for each sample in $Y$. The distribution of $Z_i$ is the distribution of possible

11

values that $r_i + \gamma \max_{a'} \tilde{Q}(s'_i, a')$ can take. From the Markov property we have that $Z_1, \ldots Z_{\lfloor \frac{|u(s,a)|}{k_m} \rfloor}$ are independent random variables.[6] From Definition 5.1 we have that $var[Z_i] \leq \sigma^2 \ \forall \ i$, and $var[G^\pi(\tilde{Q}, u(s,a), j)] \leq \frac{\sigma^2 k_m}{|u(s,a)|}$.

From Cantelli's inequality we have

$$
P\left(G^\pi(\tilde{Q}, u(s,a), j) - B^\pi \tilde{Q}(s,a) \leq -\frac{\sigma\sqrt{4k_m}}{\sqrt{|u(s,a)|}}\right)
$$

$$
\leq P\left(G^\pi(\tilde{Q}, u(s,a), j) - \mathrm{E}\left[G^\pi(\tilde{Q}, u(s,a), j)\right] \leq -\frac{\sigma\sqrt{4k_m}}{\sqrt{|u(s,a)|}}\right)
$$

$$
\leq \frac{\frac{\sigma^2 k_m}{|u(s,a)|}}{\frac{\sigma^2 k_m}{|u(s,a)|} + \left(\frac{\sigma\sqrt{4k_m}}{\sqrt{|u(s,a)|}}\right)^2}
$$

$$
= \frac{\frac{\sigma^2 k_m}{|u(s,a)|}}{\frac{\sigma^2 k_m}{|u(s,a)|} + \frac{4\sigma^2 k_m}{|u(s,a)|}}
$$

$$
= \frac{1}{5},
$$

and

$$
P\left(G^\pi(\tilde{Q}, u(s,a), j) - B^\pi \tilde{Q}(s,a) \geq \frac{\sigma\sqrt{4k_m}}{\sqrt{|u(s,a)|}}\right)
$$

$$
\leq P\left(G^\pi(\tilde{Q}, u(s,a), j) - \mathrm{E}\left[G^\pi(\tilde{Q}, u(s,a), j)\right] \geq \frac{\sigma\sqrt{4k_m}}{\sqrt{|u(s,a)|}}\right)
$$

$$
\leq \frac{\frac{\sigma^2 k_m}{|u(s,a)|}}{\frac{\sigma^2 k_m}{|u(s,a)|} + \left(\frac{\sigma\sqrt{4k_m}}{\sqrt{|u(s,a)|}}\right)^2}
$$

$$
= \frac{\frac{\sigma^2 k_m}{|u(s,a)|}}{\frac{\sigma^2 k_m}{|u(s,a)|} + \frac{4\sigma^2 k_m}{|u(s,a)|}}
$$

$$
= \frac{1}{5}.
$$

$\square$

Based on Lemma 9.4 we can now bound the probability that an individual state-action will have Bellman error of unacceptably high magnitude for a particular $\tilde{Q}$:

**Lemma 9.5.** *Let $\sigma$ be defined as in Definition 5.1, and $\epsilon_u = \max\{0, \sigma\sqrt{4k_m} - \epsilon_b\}$. For a fixed $\tilde{Q}$*

$$
P\left(F^{\pi^*}(\tilde{Q}, u(s,a)) - B^{\pi^*}\tilde{Q}(s,a) \leq -\epsilon_u\right) \leq e^{-\frac{9k_m}{50}},
$$

*and*

$$
P\left(F^{\pi^{\bar{Q}}}(\tilde{Q}, u(s,a)) - B^{\pi^{\bar{Q}}}\tilde{Q}(s,a) \geq \epsilon_u + 2\frac{\epsilon_b}{\sqrt{|u(s,a)|}}\right) \leq e^{-\frac{9k_m}{50}}.
$$

*Proof.* Let $Y$ be the set of $|u(s,a)|$ samples used by $F^\pi(\tilde{Q}, u(s,a))$ at $(s,a)$. Define $Z_1, \ldots Z_{|u(s,a)|}$ to be random variables, one for each sample in $Y$. The distribution of $Z_i$ is the distribution of next

---

[6]The state-actions the samples originate from as well as $\tilde{Q}$ and the transition model of the MDP are fixed with respect to $Z_i$, and no assumptions are made about their distribution. The only source of randomness is the the transition model of the MDP.

states $s_i'$, given $(s, a)$. From the Markov property, we have that $Z_1, \ldots Z_{|u(s,a)|}$ are independent random variables (similarly to Lemma 9.4). Let $x_j$ be a realization of $X_j$, where $X_j$'s distribution is the joint distribution of all $Z_i$ corresponding to samples that participate in $G^\pi(\tilde{Q}, u(s, a), j)$.

We define $f^{\pi^*}(x_1, \ldots x_{k_m})$ to be the function that counts the number of $j$'s such that

$$G^{\pi^*}(\tilde{Q}, u(s, a), j) - B^{\pi^*}\tilde{Q}(s, a) \leq -\frac{\sigma\sqrt{4k_m}}{\sqrt{|u(s, a)|}},$$

and $f^{\pi^{\tilde{Q}}}(x_1, \ldots x_{k_m})$ to be the function that counts the number of $j$'s such that

$$G^{\pi^{\tilde{Q}}}(\tilde{Q}, u(s, a), j) - B^{\pi^{\tilde{Q}}}\tilde{Q}(s, a) \geq \frac{\sigma\sqrt{4k_m}}{\sqrt{|u(s, a)|}}.$$

From Lemma 9.4 we have that

$$E[f^{\pi^*}(x_1, \ldots x_{k_m})] \leq \frac{k_m}{5},$$

and

$$E[f^{\pi^{\tilde{Q}}}(x_1, \ldots x_{k_m})] \leq \frac{k_m}{5}.$$

$\forall\, i \in [1, k_m]$:

$$\sup_{x_1, \ldots x_k, \hat{x}_i} |f^{\pi^*}(x_1, \ldots x_{k_a}) - f^{\pi^*}(x_1, \ldots, x_{i-1}\hat{x}_i, x_{i+1} \ldots x_{|u(s,a)|})| \leq 1,$$

and

$$\sup_{x_1, \ldots x_k, \hat{x}_i} |f^{\pi^{\tilde{Q}}}(x_1, \ldots x_{k_a}) - f^{\pi^{\tilde{Q}}}(x_1, \ldots, x_{i-1}\hat{x}_i, x_{i+1} \ldots x_{|u(s,a)|})| \leq 1.$$

From McDiarmid's inequality we have

$$P\left(f^{\pi^*}(x_1, \ldots x_{k_m}) \geq \frac{k_m}{2}\right) \leq P\left(f^{\pi^*}(x_1, \ldots x_{k_m}) - E[f^{\pi^*}(x_1, \ldots x_{k_m})] \geq \frac{3k_m}{10}\right)$$

$$\leq e^{-\frac{2\left(\frac{3k_m}{10}\right)^2}{k_m}}$$

$$= e^{-\frac{9k_m}{50}},$$

and

$$P\left(f^{\pi^{\tilde{Q}}}(x_1, \ldots x_{k_m}) \geq \frac{k_m}{2}\right) \leq P\left(f^{\pi^{\tilde{Q}}}(x_1, \ldots x_{k_m}) - E[f^{\pi^{\tilde{Q}}}(x_1, \ldots x_{k_m})] \geq \frac{3k_m}{10}\right)$$

$$\leq e^{-\frac{2\left(\frac{3k_m}{10}\right)^2}{k_m}}$$

$$= e^{-\frac{9k_m}{50}}.$$

Since the probability that

$$G^{\pi^*}(\tilde{Q}, u(s, a), j) - B^{\pi^*}\tilde{Q}(s, a) \leq -\frac{\sigma\sqrt{4k_m}}{\sqrt{|u(s, a)|}}$$

for at least $\frac{k_m}{2}$ $j$'s is bounded above by $e^{-\frac{9k_m}{50}}$, and the probability that

$$G^{\pi^{\tilde{Q}}}(\tilde{Q}, u(s, a), j) - B^{\pi^{\tilde{Q}}}\tilde{Q}(s, a) \geq \frac{\sigma\sqrt{4k_m}}{\sqrt{|u(s, a)|}}$$

for at least $\frac{k_m}{2}$ $j$'s is bounded above by $e^{-\frac{9k_m}{50}}$, the result follows from Definition 9.1. □

Given a bound on the probability that an individual state-action has Bellman error of unacceptably high magnitude, lemma 9.6 uses the union bound to bound the probability that there exists at least one state-action for some $\tilde{Q}$ produced by Median-PAC during execution, with Bellman error of unacceptably high magnitude.

**Lemma 9.6.** *Let $\epsilon_u = \max\{0, \sigma\sqrt{4k_m} - \epsilon_b\}$. The probability that for any $\tilde{Q}$ during an execution of Median-PAC there exists at least one $(s, a)$ with $|u(s, a)| > 0$ such that*

$$F^{\pi^*}(\tilde{Q}, u(s, a)) - B^{\pi^*}\tilde{Q}(s, a) \leq -\epsilon_u \tag{5}$$

*or*

$$F^{\pi^{\tilde{Q}}}(\tilde{Q}, u(s, a)) - B^{\pi^{\tilde{Q}}}\tilde{Q}(s, a) \geq \epsilon_u + 2\frac{\epsilon_b}{\sqrt{|u(s, a)|}} \tag{6}$$

*is bounded above by $2\log_2 \frac{4k}{k_m}|SA|^2 e^{-\frac{9k_m}{50}}$.*

*Proof.* At most $\log_2 \frac{4k}{k_m}|SA|$ distinct $\tilde{Q}$ exist for which $|u(s, a)| > 0$ for at least one $(s, a)$. Thus, there are at most $2\log_2 \frac{4k}{k_m}|SA|^2$ ways for at least one of the at most $|SA|$ state-actions to fail at least once during non-delay steps ($\log_2 \frac{4k}{k_m}|SA|^2$ ways each for equation 5 or equation 6 to be true at least once), each with a probability at most $e^{-\frac{9k_m}{50}}$. From the union bound, we have that the probability that for any $\tilde{Q}$ there exists at least one $(s, a)$ such that equation 5 or 6 is true, is bounded above by $2\log_2 \frac{4k}{k_m}|SA|^2 e^{-\frac{9k_m}{50}}$. $\square$

Based on Lemma 9.6 we can now bound the probability that any $(s, a)$ will have Bellman error of unacceptably high magnitude:

**Lemma 9.7.** *Let $\epsilon_u = \max\{0, \sigma\sqrt{4k_m} - \epsilon_b\}$. The probability that for any $\tilde{Q}$ during an execution of Median-PAC there exists at least one $(s, a)$ such that*

$$\tilde{Q}(s, a) - B^{\pi^*}\tilde{Q}(s, a) \leq -\epsilon_u - \epsilon_a \tag{7}$$

*or at least one $(s, a)$ with $|u(s, a)| > 0$ such that*

$$\tilde{Q}(s, a) - B^{\pi^{\tilde{Q}}}\tilde{Q}(s, a) \geq \epsilon_u + \epsilon_a + 2\frac{\epsilon_b}{\sqrt{|u(s, a)|}} \tag{8}$$

*is bounded above by $2\log_2 \frac{4k}{k_m}|SA|^2 e^{-\frac{9k_m}{50}}$.*

*Proof.* When $|u(s, a)| < k_m$, $\tilde{Q}(s, a) = Q_{\max}$. Since $B^{\pi^*}\tilde{Q}(s, a) \leq Q_{\max}$, $\tilde{Q}(s, a) - B^{\pi^*}\tilde{Q}(s, a) \leq -\epsilon_u - \epsilon_a$. Otherwise, $\forall(s, a, \tilde{Q})$ with probability $1 - 2\log_2 \frac{4k}{k_m}|SA|^2 e^{-\frac{9k_m}{50}}$

$$\begin{aligned} B^{\pi^*}\tilde{Q}(s, a) &= \min\left\{Q_{\max}, B^{\pi^*}\tilde{Q}(s, a)\right\} \\ &< \min\left\{Q_{\max}, F^{\pi^*}(\tilde{Q}, u(s, a)) + \epsilon_u\right\} \\ &\leq \tilde{B}^{\pi^*}\tilde{Q}(s, a) + \epsilon_u \\ &\leq \tilde{B}\tilde{Q}(s, a) + \epsilon_u \\ &\leq \tilde{Q}(s, a) + \epsilon_u + \epsilon_a. \end{aligned}$$

14

$\forall (s,a,\tilde{Q})$ with $u(s,a) \geq k_m$, with probability $1 - 2\log_2 \frac{4k}{k_m}|SA|^2 e^{-\frac{9k_m}{50}}$

$$
\begin{aligned}
B^{\pi^{\bar{Q}}}\tilde{Q}(s,a) &= \min\left\{Q_{\max}, B^{\pi^{\bar{Q}}}\tilde{Q}(s,a)\right\} \\
&> \min\left\{Q_{\max}, F^{\pi^{\bar{Q}}}(\tilde{Q}, u(s,a)) - \epsilon_u - 2\frac{\epsilon_b}{\sqrt{|u(s,a)|}}\right\} \\
&\geq \min\left\{Q_{\max}, F^{\pi^{\bar{Q}}}(\tilde{Q}, u(s,a))\right\} - \epsilon_u - 2\frac{\epsilon_b}{\sqrt{|u(s,a)|}} \\
&\geq \tilde{B}^{\pi^{\bar{Q}}}\tilde{Q}(s,a) - \epsilon_u - 2\frac{\epsilon_b}{\sqrt{|u(s,a)|}} \\
&= \tilde{B}\tilde{Q}(s,a) - \epsilon_u - 2\frac{\epsilon_b}{\sqrt{|u(s,a)|}} \\
&\geq \tilde{Q}(s,a) - \epsilon_u - \epsilon_a - 2\frac{\epsilon_b}{\sqrt{|u(s,a)|}}.
\end{aligned}
$$

Note that both the first half of Lemma 9.6 (used in the fist half of the proof) and the second half (used in the second half of the proof) hold simultaneously with probability $2\log_2 \frac{4k}{k_m}|SA|^2 e^{-\frac{9k_m}{50}}$, therefore we do not need to take a union bound over the individual probabilities. $\qquad\square$

We will use the following three lemmas from Pazis and Parr (2016):

**Lemma 9.8.** *Let $t_i$ for $i = 0 \to l$ be the outcomes of independent (but not necessarily identically distributed) random variables in $\{0,1\}$, with $P(t_i = 1) \geq p_i$. If $\frac{2}{m}\ln\frac{1}{\delta} < 1$ and*

$$
\sum_{i=0}^{l} p_i \geq \frac{m}{1 - \sqrt{\frac{2}{m}\ln\frac{1}{\delta}}},
$$

*then $\sum_{i=0}^{l} t_i \geq m$ with probability at least $1 - \delta$.*

**Lemma 9.9.** *Let $Q(s,a) - B^{\pi^*}Q(s,a) \geq -\epsilon_* \; \forall(s,a)$, $X_1, \ldots, X_i, \ldots, X_n$ be sets of state-actions where $Q(s,a) - B^{\pi^Q}Q(s,a) \leq \epsilon_i \; \forall(s,a) \in X_i$, $Q(s,a) - B^{\pi^Q}Q(s,a) \leq \epsilon_{\pi Q} \; \forall(s,a) \notin \cup_{i=1}^{n}X_i$, and $\epsilon_{\pi Q} \leq \epsilon_i \forall i$. Let $T_H = \left\lceil \frac{1}{1-\gamma}\ln\frac{(1-\gamma)Q_{\max}}{\epsilon_a}\right\rceil$ and define $H = \{1,2,4,\ldots,2^i\}$ where $i$ is the largest integer such that $2^i \leq T_H$. Define $p_{h,i}(s)$ for $h \in [0, T_H-1]$ to be Bernoulli random variables expressing the probability of encountering exactly $h$ state-actions for which $(s,a) \in X_i$ when starting from state $s$ and following $\pi^Q$ for a total of $\min\{T, T_H\}$ steps. Finally let $p_{h,i}^e(s) = \sum_{m=h}^{2h-1} p_{m,i}(s)$. Then*

$$
V^*(s) - V^{\pi^Q}(s) \leq \frac{\epsilon_* + \epsilon_{\pi Q} + \epsilon_a}{1-\gamma} + \epsilon_e,
$$

*where $\epsilon_e = 2\sum_{i=1}^{n}\left(\sum_{h \in H}\left(hp_{h,i}^e(s)\right)(\epsilon_i - \epsilon_{\pi Q})\right) + \gamma^T Q_{\max}$.*

**Lemma 9.10.** *Let $\hat{B}$ be a $\gamma$-contraction with fixed point $\hat{Q}$, and $Q$ the output of*

$$
\frac{1}{1-\gamma}\ln\frac{Q_{\max}}{\epsilon}
$$

*iterations of value iteration using $\hat{B}$. Then if $0 \leq \hat{Q}(s,a) \leq Q_{\max}$ and $0 \leq Q_0(s,a) \leq Q_{\max} \; \forall(s,a)$, where $Q_0(s,a)$ is the initial value for $(s,a)$*

$$
-\epsilon \leq Q(s,a) - \hat{B}Q(s,a) \leq \epsilon \; \forall(s,a).
$$

Lemma 9.11 bounds the number of times the policy produced by Median-PAC can encounter state-actions with fewer than $k$ samples.

**Lemma 9.11.** *Let $(s_1, s_2, s_3, \ldots)$ be the random path generated on some execution of Algorithm 1. Let $\tau(t)$ be the number of steps from step $t$ to the next step for which the policy changes. Let $T_H = \left\lceil \frac{1}{1-\gamma}\ln\frac{(1-\gamma)Q_{\max}}{\epsilon_a}\right\rceil$ and define $H = \{1,2,4,\ldots,2^i\}$ where $i$ is the largest such*

that $2^i \leq T_H$. Let $K_a = \{2^0 k_m, 2^1 k_m, 2^2 k_m, \ldots k\}$. Let $k_a^-$ be the largest value in $K_a$ that is strictly smaller than $k_a$, or 0 if such a value does not exist. Let $X_{k_a}(t)$ be the set of state-actions at step $t$ for which $k_a^- = |u(s,a)|$. Define $p_{h,k_a}(s_t)$ for $k_a \in K_a$ to be Bernoulli random variables that express the following conditional probability: Given $\tilde{Q}$ at step $t$, exactly $h$ state-actions in $X_{k_a}(t)$ are encountered during the next $\min\{T_H, \tau(t)\}$ steps. Let $p_{h,k_a}^e(s_t) = \sum_{i=h}^{2h-1} p_{i,k_a}(s_t)$. If

$$\frac{2\left\lceil \frac{1}{1-\gamma} \ln \frac{(1-\gamma)Q_{\max}}{\epsilon_a} \right\rceil 2 \ln \frac{\log_2 \frac{2k}{k_m}}{\delta}}{k_m|SA|+1} < 1, \text{ with probability at least } 1-\delta$$

$$\sum_{t=0}^{\infty} \sum_{h \in H} (h p_{h,k_a}^e(s_{t,j})) < \frac{(k_a|SA|+1)\left(1 + \log_2 \left\lceil \frac{1}{1-\gamma} \ln \frac{(1-\gamma)Q_{\max}}{\epsilon_a} \right\rceil\right) \left\lceil \frac{1}{1-\gamma} \ln \frac{(1-\gamma)Q_{\max}}{\epsilon_a} \right\rceil}{1 - \sqrt{\frac{2\left\lceil \frac{1}{1-\gamma} \ln \frac{(1-\gamma)Q_{\max}}{\epsilon_a} \right\rceil 2 \ln \frac{\log_2 \frac{2k}{k_m}}{\delta}}{k_m|SA|+1}}}$$

$\forall k_a \in K_a$ and $\forall h \in H$ simultaneously.

*Proof.* From the Markov property we have that $p_{h,k_a}^e(s_t)$ variables at least $T_H$ steps apart are independent.[7] Define $T_i^H$ for $i \in \{0, 1, \ldots, T_H - 1\}$ to be the (infinite) set of timesteps for which $t \in \{i, i + T_H, i + 2T_H, \ldots\}$.

Since $k_a$ samples will be added to a state-action such that $|u(s,a)| = k_a^-$ before $|u(s,a)| = k_a$, at most $k_a|SA|$ state-actions such that $k_a^- = |u(s,a)|$ can be encountered.

Let us assume that there exists an $i \in \{0, 1, \ldots, T_H - 1\}$ and $h \in H$ such that

$$\sum_{t \in T_i^H} p_{h,k_a}^e(s_{t,j}) \geq \frac{k_a|SA|+1}{h\left(1 - \sqrt{\frac{2h}{k_a|SA|+1} \ln \frac{2\log_2 \frac{2k}{k_m}}{\delta}}\right)}.$$

From Lemma 9.8 it follows that with probability at least $1 - \frac{\delta}{2\log_2 \frac{2k}{k_m}}$, at least $k_a|SA|+1$ state-actions such that $k_a^- = |u(s,a)|$ will be encountered, which is a contradiction. It must therefore be the case that

$$\sum_{t \in T_i^H} p_{h,k_a}^e(s_{t,j}) < \frac{k_a|SA|+1}{h\left(1 - \sqrt{\frac{2h}{k_a|SA|+1} \ln \frac{2\log_2 \frac{2k}{k_m}}{\delta}}\right)}$$

with probability at least $1 - \frac{\delta}{2\log_2 \frac{2k}{k_m}}$ for all $i \in \{0, 1, \ldots, T_H - 1\}$ and $h \in H - \{T_H\}$ simultaneously, which implies that

$$\sum_{t=0}^{\infty} \sum_{h \in H} (h p_{h,k_a}^e(s_{t,j})) < \frac{(k_a|SA|+1)|H|T_H}{1 - \sqrt{\frac{2T_H}{k_a|SA|+1} \ln \frac{2\log_2 \frac{2k}{k_m}}{\delta}}}$$

$$\leq \frac{(k_a|SA|+1)\left(1 + \log_2 \left\lceil \frac{1}{1-\gamma} \ln \frac{(1-\gamma)Q_{\max}}{\epsilon_a} \right\rceil\right) \left\lceil \frac{1}{1-\gamma} \ln \frac{(1-\gamma)Q_{\max}}{\epsilon_a} \right\rceil}{1 - \sqrt{\frac{2\left\lceil \frac{1}{1-\gamma} \ln \frac{(1-\gamma)Q_{\max}}{\epsilon_a} \right\rceil 2 \ln \frac{\log_2 \frac{2k}{k_m}}{\delta}}{k_m|SA|+1}}}$$

with probability at least $1 - \frac{\delta}{2\log_2 \frac{2k}{k_m}}$ for all $h \in H$ simultaneously.

From the union bound we have that since $k_a$ can take at most $\log_2 \frac{2k}{k_m}$ values, with probability $1 - \delta$

$$\sum_{t=0}^{\infty} \sum_{h \in H} (h p_{h,k_a}^e(s_{t,j})) < \frac{(k_a|SA|+1)\left(1 + \log_2 \left\lceil \frac{1}{1-\gamma} \ln \frac{(1-\gamma)Q_{\max}}{\epsilon_a} \right\rceil\right) \left\lceil \frac{1}{1-\gamma} \ln \frac{(1-\gamma)Q_{\max}}{\epsilon_a} \right\rceil}{1 - \sqrt{\frac{2\left\lceil \frac{1}{1-\gamma} \ln \frac{(1-\gamma)Q_{\max}}{\epsilon_a} \right\rceil 2 \ln \frac{\log_2 \frac{2k}{k_m}}{\delta}}{k_m|SA|+1}}}$$

$\forall k_a \in K_a$ and $\forall h \in H$ simultaneously. $\square$

---

[7]While what happens at step $t$ affects which variables are selected at future timesteps, this is not a problem. We only care that the *outcomes* of the variables are independent given their selection.

**Lemma 5.3.** *The per step computational complexity of algorithm 1 is bounded above by:*

$$O\left(\frac{k|S||A|^2}{1-\gamma}\ln\frac{Q_{\max}}{\epsilon_a}\right).$$

*Proof.* From lemma 9.10 we have that on every iteration of algorithm 1, lines 13 through 15 will we executed at most $O\left(\frac{1}{1-\gamma}\ln\frac{Q_{\max}}{\epsilon_a}\right)$ times. For each one of these iterations, function $\tilde{B}\tilde{Q}(s,a)$ will be called $|S||A|$ times. Line 22 in function $\tilde{B}\tilde{Q}(s,a)$ will be executed at most $k_m$ times, with a per execution cost of $O\left(\frac{k}{k_m}|A|\right)$. $\square$

**Theorem 5.4.** *Let $(s_1, s_2, s_3, \dots)$ be the random path generated on some execution of Median-PAC, and $\tilde{\pi}$ be the (non-stationary) policy followed by Median-PAC. Let $\epsilon_u = \max\{0, \sigma\sqrt{4k_m} - \epsilon_a\sqrt{k}\}$, and $\epsilon_a$ be defined as in algorithm 1. If $k_m = \left\lceil \frac{50}{9}\ln\frac{4\log_2\frac{4Q_{\max}^2}{\epsilon_a^2}|SA|^2}{\delta}\right\rceil$,*

*$\frac{2\left\lceil\frac{1}{1-\gamma}\ln\frac{(1-\gamma)Q_{\max}}{\epsilon_a}\right\rceil 2\ln\frac{\log_2\frac{2k}{k_m}}{\delta}}{k_m|SA|+1} < 1$, and $k = 2^i k_m$ for some integer $i$, then with probability at least $1-\delta$, for all $t$*

$$V^*(s_t) - V^{\tilde{\pi}}(s_t) \leq \frac{2\epsilon_u + 5\epsilon_a}{1-\gamma} + \epsilon_e(t),$$

*where*

$$\sum_{t=0}^{\infty}\epsilon_e(t) < c_0\left(\left(2k_m + \log_2\frac{2k}{k_m}\right)Q_{\max} + \epsilon_a k\left(8 + \frac{8}{\sqrt{2}}\right)\right),$$

*and*

$$c_0 = \frac{(|SA|+1)\left(1 + \log_2\left\lceil\frac{1}{1-\gamma}\ln\frac{(1-\gamma)Q_{\max}}{\epsilon_a}\right\rceil\right)\left\lceil\frac{1}{1-\gamma}\ln\frac{(1-\gamma)Q_{\max}}{\epsilon_a}\right\rceil}{1 - \sqrt{\frac{2\left\lceil\frac{1}{1-\gamma}\ln\frac{(1-\gamma)Q_{\max}}{\epsilon_a}\right\rceil 2\ln\frac{\log_2\frac{2k}{k_m}}{\delta}}{k_m|SA|+1}}}.$$

*If $k = 2^i k_m$ where $i$ is the smallest integer such that $2^i \geq \frac{4\sigma^2}{\epsilon_a^2}$, and $\epsilon_0 = (1-\gamma)\epsilon_a$, then with probability at least $1-\delta$, for all $t$*

$$V^*(s_t) - V^{\tilde{\pi}}(s_t) \leq \epsilon_0 + \epsilon_e(t),$$

*where[8]*

$$\sum_{t=0}^{\infty}\epsilon_e(t) \approx \tilde{O}\left(\left(\frac{\sigma^2}{\epsilon_0(1-\gamma)^2} + \frac{Q_{\max}}{1-\gamma}\right)|SA|\right).$$

*Note that the probability of success holds for all timesteps simultaneously, and $\sum_{t=0}^{\infty}\epsilon_e(t)$ is an undiscounted infinite sum.*

*Proof.* From Lemma 9.7 we have that with probability at least $1 - 2\log_2\frac{4k}{k_m}|SA|^2 e^{-\frac{9k_m}{50}}$

$$\tilde{Q}(s,a) - B^{\pi^*}\tilde{Q}(s,a) > -\epsilon_u - \epsilon_a \tag{9}$$

for all $(s, a, \tilde{Q})$, and

$$\tilde{Q}(s,a) - B^{\pi^{\tilde{Q}}}\tilde{Q}(s,a) < \epsilon_u + \epsilon_a + 2\frac{\epsilon_b}{\sqrt{|u(s,a)|}} \tag{10}$$

for all $(s, a)$ with $|u(s,a)| \geq k_m$. We also have that

$$\tilde{Q}(s,a) - B^{\pi^{\tilde{Q}}}\tilde{Q}(s,a) \leq Q_{\max} - Q_{\min} \ \forall\ (s,a,\tilde{Q}).$$

---

[8] $f(n) = \tilde{O}(g(n))$ is a shorthand for $f(n) = O(g(n)\log^c g(n))$ for some $c$.

Let $K_a$, $k_a^-$, $T_H$, $H$, $\tau(t)$, and $p_{h,k_a}^e(s_t)$ be defined as in lemma 9.11. With probability at least $1 - 2\log_2 \frac{4k}{k_m}|SA|^2 e^{-\frac{9k_m}{50}}$, for any $(s,a)$ with $|u(s,a)| > 0$ samples

$$\tilde{Q}(s,a) - B^{\pi^{\tilde{Q}}}\tilde{Q}(s,a) < \epsilon_u + \epsilon_a + 2\frac{\epsilon_b}{\sqrt{|u(s,a)|}}.$$

Even though $\tilde{\pi}$ is non-stationary, it is comprised of stationary segments. Starting from step $t$, $\tilde{\pi}$ is stationary for at least $\tau(t)$ steps. Substituting the above into Lemma 9.9 we have that with probability at least $1 - 2\log_2 \frac{4k}{k_m}|SA|^2 e^{-\frac{9k_m}{50}}$

$$V^*(s_t) - V^{\tilde{\pi}}(s_t) \leq \frac{2\epsilon_u + 3\epsilon_a + 2\frac{\epsilon_b}{\sqrt{k}}}{1-\gamma} + \epsilon_e(t),$$

where

$$\epsilon_e(t) = \gamma^{\tau(t)}Q_{\max} + 2\sum_{h\in H}(hp_{h,k_m}^e(s_t))Q_{\max} + \sum_{k_a\in\{K_a-k_m\}}2\sum_{h\in H}(hp_{h,k_a}^e(s_t))2\frac{\epsilon_b}{\sqrt{k_a}}.$$

From the above it follows that

$$\sum_{t=0}^{\infty}\epsilon_e(t)$$

$$= \sum_{t=0}^{\infty}\left(\gamma^{\tau(t)}Q_{\max} + 2\sum_{h\in H}(hp_{h,1}^e(s_{t,j}))Q_{\max} + \sum_{k_a\in\{K_a-k_m\}}2\sum_{h\in H}(hp_{h,k_a}^e(s_{t,j}))2\frac{\epsilon_b}{\sqrt{k_a}}\right)$$

$$= \sum_{t=0}^{\infty}\gamma^{\tau(t)}Q_{\max} + 2\sum_{t=0}^{\infty}\sum_{h\in H}(hp_{h,1}^e(s_{t,j}))Q_{\max} + 2\sum_{k_a\in\{K_a-k_m\}}\sum_{t=0}^{\infty}\sum_{h\in H}(hp_{h,k_a}^e(s_{t,j}))2\frac{\epsilon_b}{\sqrt{k_a}}$$

$$< \frac{|SA|Q_{\max}\log_2\frac{2k}{k_m}}{(1-\gamma)} + 2k_m c_0 Q_{\max} + 2\sum_{k_a\in\{K_a-k_m\}}k_a c_0 2\frac{\epsilon_b}{\sqrt{k_a}}$$

$$< \left(2k_m + \log_2\frac{2k}{k_m}\right)c_0 Q_{\max} + 2\sum_{k_a\in\{K_a-k_m\}}k_a c_0 2\frac{\epsilon_b}{\sqrt{k_a}}$$

$$= \left(2k_m + \log_2\frac{2k}{k_m}\right)c_0 Q_{\max} + 4c_0\epsilon_b\sum_{k_a\in\{K_a-k_m\}}\sqrt{k_a}$$

$$< \left(2k_m + \log_2\frac{2k}{k_m}\right)c_0 Q_{\max} + 4c_0\epsilon_b\sqrt{k}\left(\sum_{i=0}^{\infty}\left(\frac{1}{2^i} + \frac{1}{2^i\sqrt{2}}\right)\right)$$

$$= c_0\left(\left(2k_m + \log_2\frac{2k}{k_m}\right)Q_{\max} + \epsilon_b\sqrt{k}\left(8 + \frac{8}{\sqrt{2}}\right)\right)$$

with probability $1 - \delta - 2\log_2\frac{4k}{k_m}|SA|^2 e^{-\frac{9k_m}{50}}$, where in step 3 we used the fact that there can be at most $\log_2\frac{2k}{k_m}|SA|$ policy changes. Since Lemma 9.7 (used to bound the Bellman error of each $(s,a,\tilde{Q})$) holds with probability $2\log_2\frac{4k}{k_m}|SA|^2 e^{-\frac{9k_m}{50}}$ and Lemma 9.11 (used to bound how many times each $(s,a,\tilde{Q})$ is encountered) holds with probability of at least $1 - \delta$, the bound above holds with probability of at least $1 - \delta - 2\log_2\frac{4k}{k_m}|SA|^2 e^{-\frac{9k_m}{50}}$.

Setting $\epsilon_b = \epsilon_a\sqrt{k}$ we have that with probability at least $1 - \delta - 2\log_2\frac{4k}{k_m}|SA|^2 e^{-\frac{9k_m}{50}}$

$$V^*(s_t) - V^{\tilde{\pi}}(s_t) \leq \frac{2\epsilon_u + 5\epsilon_a}{1-\gamma} + \epsilon_e(t),$$

where

$$\sum_{t=0}^{\infty}\epsilon_e(t) < c_0\left(\left(2k_m + \log_2\frac{2k}{k_m}\right)Q_{\max} + \epsilon_a k\left(8 + \frac{8}{\sqrt{2}}\right)\right).$$

Equations 3 and 4 follow by substitution and by using the fact that $\sigma \leq \frac{Q_{\max}}{2}$. $\qquad\square$