

A Omitted Related Work

Neural Kernel Bandits. [66] initiated the study of kernelized linear bandits, showing regret dependent on the information gain. The work of [76, 73] specialized this to the Neural Tangent Kernel (NTK) [54, 24, 40, 27, 10], where the algorithm utilizes gradient descent but remains close to initialization and thus remains a kernel class. Furthermore, NTK methods require d^p samples to express a degree p polynomial in d dimensions [31], similar to eluder dimension of polynomials, and so lack the inductive biases necessary for real-world applications of decision-making problems [60].

Neural Bandits. For bandits with practical neural networks (instead of overparameterized NTKs) as the function approximator, we are not aware of any previous paper that gives provably efficient bandit algorithm for this case. Our paper gives the first provably efficient algorithm for neural bandits with noiseless reward and deterministic activation. We note that however, previous paper has already solved neural bandits when the neural network happens to be convex [21].

Concave Bandits. There has been a rich line of work on concave bandits starting with [26, 46]. [4] attained the first \sqrt{T} regret algorithm for concave bandits though with a large $\text{poly}(d)$ dependence. In the adversarial setting, a line of work [38, 14, 49] have attained polynomial-time algorithms with \sqrt{T} regret with increasingly improved dimension dependence. The sharp dimension dependence remains unknown.

Noiseless Bandits. In the noiseless setting, there is some investigation in phase retrieval borrowing the tools from algebraic geometry (see e.g. [69]). In this paper, we will study the bandit problem with more general reward functions: neural nets with polynomial activation (structured polynomials) including phase retrieval. [45] study similar structured polynomials, also using tools from algebraic geometry, but they only study the expressivity of those polynomials and do not consider the learning problems. [21] study noiseless bandits with bounded Sequential Rademacher Complexity, but focus on attaining local optimality.

Concurrent work. [50] address the phase retrieval bandit problem which is equivalent to a symmetric rank 1 variant of the bilinear bandit of [43] and attain $\tilde{O}(\sqrt{d^2 T})$ regret. Our work in Section 3.1 specialized to the rank 1 case attains the same regret.

Matrix/Tensor Power Method. Our analysis stems from noisy power methods for matrix/tensor decomposition problems. Robust power method, subspace iteration, and tensor decomposition that tolerate noise first appeared in [37, 8]. Follow-up work attained the optimal rate for both gap-dependence and gap-free settings for matrix decomposition [58, 6]. An improvement on the problem dimension for tensor power method is established in [70]. [63] considers the convergence of tensor power method in the non-orthogonal case.

B Additional Preliminaries

In this section we show that adapting the eluder UCB algorithms from [62] would yield the sample complexity in Theorem 2.1. Especially we give the rates in Table 1 for our stochastic settings.

The algorithm [62] consider Algorithm 2 for the stochastic generalized linear bandit problem. Assume that θ^* is the true parameter of the reward model. The reward is $r_t = f_{\theta^*}(\mathbf{a}_t) + \eta_t$ for $f_{\theta^*} \in \mathcal{F}$. Let N be the α -covering-number (under $\|\cdot\|_\infty$) of \mathcal{F} , d_E be the α -eluder-dimension of \mathcal{F} (see Definition 3,4 in [62]). Let $C = \sup_{f \in \mathcal{F}, a \in \mathcal{A}} |f(a)|$. We set $\alpha = \frac{1}{T^2}$ in the algorithm.

The regret analysis Choosing $\alpha = 1/T^2$, proposition 4 in [62] state that with probability $1 - \delta$, for some universal constant C , the total regret $\mathfrak{R}(T) \leq \frac{1}{T} + C \min\{d_E, T\} + 4\sqrt{d_E \beta_T T} \leq 1 + C\sqrt{d_E T} + 4\sqrt{d_E \beta_T T} = O(\sqrt{d_E(1 + \beta_T T)})$. In our settings with $\alpha = 1/T^2$, $\beta_T = 8 \log(N/\delta) + 2(8C + \sqrt{8 \ln(4T^2/\delta)})/T = O(\log(N/\delta))$ where $\log(N) = \Omega(1)$ for our action sets, and thus

$$\mathfrak{R}(T) = \tilde{O}(\sqrt{d_E T \log N}).$$

Algorithm 2 Eluder UCB

- 1: **Input:** Function class \mathcal{F} , failure probability δ , parameters α, N, C .
 - 2: **Initialization:** $\mathcal{F}_0 \leftarrow \mathcal{F}$.
 - 3: **for** t from 1 to T **do**
 - 4: **Select Action:**
 - 5: $\mathbf{a}_t \in \arg \max_{\mathbf{a} \in \mathcal{A}} \sup_{f_{\theta} \in \mathcal{F}_{t-1}} f_{\theta}(\mathbf{a})$
 - 6: Play action \mathbf{a}_t and observe reward r_t
 - 7: **Update Statistics:**
 - 8: $\widehat{\theta}_t \in \arg \min_{\theta} \sum_{s=1}^t (f_{\theta}(\mathbf{a}_s) - r_s)^2$
 - 9: $\beta_t \leftarrow 8 \log(N/\delta) + 2\alpha t(8C + \sqrt{8 \ln(4t^2/\delta)})$
 - 10: $\mathcal{F}_t \leftarrow \{f_{\theta} : \sum_{s=1}^t (f_{\theta} - f_{\widehat{\theta}_t})^2(\mathbf{a}_s) \leq \beta_t\}$
-

Applications in our settings We show that in our settings Theorem 2.1 will obtain the rates listed in Table 1.

The covering numbers

Lemma B.1. *The log-covering-number (of radius α with $\alpha \ll 1$, under $\|\cdot\|_{\infty}$) of the function classes are: $\log N(\mathcal{F}_{\text{SYM}}) = O(dk \log \frac{k}{\alpha})$, $\log N(\mathcal{F}_{\text{ASYM}}) = O(dk \log \frac{k}{\alpha})$, $\log N(\mathcal{F}_{\text{EV}}) = O(dk \log \frac{k}{\alpha})$, and $\log N(\mathcal{F}_{\text{LR}}) = O(dk \log \frac{k}{\alpha})$.*

Proof. Let S_{ξ}^d denote a minimal ξ -covering of \mathbb{S}^{d-1} (under $\|\cdot\|_2$) for $0 < \xi < \frac{1}{10}$, and $|S_{\xi}^d| = O(d \log 1/\xi)$ (see for example [62]). Then we can construct the coverings in our settings from S_{ξ}^d :

- \mathcal{F}_{SYM} : let $\xi = \frac{\alpha}{kp}$, and for k copies of S_{ξ}^d , we can construct a covering of \mathcal{F}_{SYM} with size $|S_{\xi}^d|^k$. Specifically, let the covering be $S_{\text{SYM}} = \{g(\mathbf{a}) = \sum_{j=1}^k \lambda_j (\mathbf{u}_j^{\top} \mathbf{a})^p : (\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k) \in S_{\xi}^d \times S_{\xi}^d \times \dots \times S_{\xi}^d\}$, then for each $f(\mathbf{a}) = \sum_{j=1}^k \lambda_j (\mathbf{v}_j^{\top} \mathbf{a})^p \in \mathcal{F}_{\text{SYM}}$, as we can find $\mathbf{u}_j \in S_{\xi}^d$ that $\|\mathbf{u}_j - \mathbf{v}_j\|_2 \leq \xi$,

$$\sup_{\mathbf{a}} [f(\mathbf{a}) - g(\mathbf{a})] \leq \sup_{\mathbf{a}} \left[\sum_{j=1}^k |\lambda_j| \|\mathbf{u}_j^{\top} \mathbf{a} - \mathbf{v}_j^{\top} \mathbf{a}\| \sum_{q=0}^{p-1} (\mathbf{u}_j^{\top} \mathbf{a})^q (\mathbf{v}_j^{\top} \mathbf{a})^{p-q-1} \right] \leq pk\xi = \alpha;$$

- $\mathcal{F}_{\text{ASYM}}$: let $\xi = \frac{\alpha}{kp}$, and for kp copies of S_{ξ}^d , let the covering be $S_{\text{ASYM}} = \{g(\mathbf{a}) = \sum_{j=1}^k \lambda_j \prod_{q=1}^p (\mathbf{u}_j(q)^{\top} \mathbf{a}(q)) : (\mathbf{u}_1(1), \mathbf{u}_1(2), \dots, \mathbf{u}_1(p), \mathbf{u}_2(1), \dots, \mathbf{u}_k(p)) \in S_{\xi}^d \times S_{\xi}^d \times \dots \times S_{\xi}^d\}$ with size $|S_{\xi}^d|^{kp}$. Then for each $f(\mathbf{a}) = \sum_{j=1}^k \lambda_j \prod_{q=1}^p (\mathbf{v}_j(q)^{\top} \mathbf{a}(q)) \in \mathcal{F}_{\text{ASYM}}$, as we can find $\mathbf{u}_j(q) \in S_{\xi}^d$ that $\|\mathbf{u}_j(q) - \mathbf{v}_j(q)\|_2 \leq \xi$,

$$\begin{aligned} \sup_{\mathbf{a}} [f(\mathbf{a}) - g(\mathbf{a})] &\leq \sup_{\mathbf{a}} \left[\sum_{j=1}^k |\lambda_j| \sum_{q=1}^p \|\mathbf{u}_j(q)^{\top} \mathbf{a} - \mathbf{v}_j(q)^{\top} \mathbf{a}\| \cdot \right. \\ &\quad \left. \left| \prod_{r < q} (\mathbf{u}_j(r)^{\top} \mathbf{a}) \prod_{r > q} (\mathbf{v}_j(r)^{\top} \mathbf{a}) \right| \right] \\ &\leq pk\xi = \alpha; \end{aligned}$$

- \mathcal{F}_{EV} : the construction follows that of \mathcal{F}_{SYM} by taking $p = 2$;
- \mathcal{F}_{LR} : taking the construction of \mathcal{F}_{SYM} with $p = 2$ and $\xi = \frac{\alpha}{2k}$, for $\mathbf{N} = \sum_{j=1}^k \lambda_j \mathbf{u}_j \mathbf{u}_j^{\top}$ and $\mathbf{M} = \sum_{j=1}^k \lambda_j \mathbf{v}_j \mathbf{v}_j^{\top}$ with $\|\mathbf{u}_j - \mathbf{v}_j\|_2 \leq \xi$, we know $\|\mathbf{N} - \mathbf{M}\|_{\text{F}} \leq \left\| \mathbf{N} - \sum_{j=1}^k \lambda_j \mathbf{u}_j \mathbf{v}_j^{\top} \right\|_{\text{F}} + \left\| \sum_{j=1}^k \lambda_j \mathbf{u}_j \mathbf{v}_j^{\top} - \mathbf{M} \right\|_{\text{F}} \leq \sum_{j=1}^k 2|\lambda_j| \xi \leq \alpha$. Then $\sup_{\mathbf{A}} [f_{\mathbf{M}}(\mathbf{A}) - f_{\mathbf{N}}(\mathbf{A})] \leq \sup_{\mathbf{A}} \|\mathbf{M} - \mathbf{N}\|_{\text{F}} \cdot \|\mathbf{A}\|_{\text{F}} \leq \alpha$.

Then we can bound the covering numbers in Theorem 2.1. Notice that in the settings the log-covering numbers are only different by constant factors. \square

The eluder dimensions

Lemma B.2. *The ϵ -eluder-dimension ($\epsilon < 1$) d_E of the function classes are: $d_E(\mathcal{F}_{\text{SYM}}) = \tilde{\Theta}(d^p)$ (for $k \geq p$), $d_E(\mathcal{F}_{\text{ASYM}}) = \tilde{\Theta}(d^p)$, $d_E(\mathcal{F}_{\text{EV}}) = \tilde{\Theta}(d^2)$, and $d_E(\mathcal{F}_{\text{LR}}) = \tilde{\Theta}(d^2)$. In the settings WLOG we assume the top eigenvalue is $r^* = \lambda_1 = 1$ as we are mostly interested in the cases where $r^* > \epsilon$.*

Proof. The upper bounds for the eluder dimension can be given by the linear argument. [62] show that the d -dimension linear model $\{f_{\boldsymbol{\theta}}(\mathbf{a}) = \boldsymbol{\theta}^\top \mathbf{a}\}$ has ϵ -eluder-dimension $O(d \log \frac{1}{\epsilon})$. In all of these settings, we can find feature maps ϕ and ψ so that $\mathcal{F} = \{f_{\boldsymbol{\theta}}(\mathbf{a}), f_{\boldsymbol{\theta}}(\mathbf{a}) = \phi(\boldsymbol{\theta})^\top \psi(\mathbf{a}), \|\phi(\boldsymbol{\theta})\|_2 \leq k, \|\psi(\mathbf{a})\|_2 \leq k\}$. Then the eluder dimensions will be bounded by the corresponding linear dimension as an original ϵ -independent sequence $\{\mathbf{a}_i\}$ will induce an ϵ -independent sequence $\{\psi(\mathbf{a}_i)\}$ in the linear model. Therefore for matrices (\mathcal{F}_{LR} and \mathcal{F}_{EV}) the eluder dimension is $O(d^2 \log \frac{k}{\epsilon})$ and for the tensors (\mathcal{F}_{SYM} and $\mathcal{F}_{\text{ASYM}}$) it is $O(d^p \log \frac{k}{\epsilon})$.

Then we consider the lower bounds. We provide the following example of $O(1)$ -independent sequences to bound the eluder dimension in our settings up to a log factor.

- \mathcal{F}_{SYM} : the sequence is $\{\mathbf{a}_i = (\mathbf{e}_{i_1}, \mathbf{e}_{i_2}, \dots, \mathbf{e}_{i_p}) : i = (i_1, i_2, \dots, i_p) \in [d]^p\}$. For $f_j(\mathbf{a}) = \prod_{q=1}^p \mathbf{e}_{j_q}^\top \mathbf{a}(q)$, $f_j(\mathbf{a}_i)$ is only 1 when $i = j$ and 0 otherwise. Then each \mathbf{a}_i is 1-independent to the predecessors on f_i and zero, and thus the eluder dimension is lower bounded by d^p .
- $\mathcal{F}_{\text{ASYM}}$: for $p \leq d$ and $k \geq p$, the sequence is $\{\mathbf{a}_i = \frac{1}{\sqrt{p}}(\mathbf{e}_{i_1} + \mathbf{e}_{i_2} + \dots + \mathbf{e}_{i_p}) : i = (i_1, i_2, \dots, i_p) \in [d]^p, i_1 < i_2 < \dots < i_p\}$. There are tensors f_j and g_j of CP-rank k that $(f_j - g_j)(\mathbf{a}) = \prod_{q=1}^p (\mathbf{e}_{j_q}^\top \mathbf{a})$ where $j_1 < j_2 < \dots < j_p$, $(f_j - g_j)(\mathbf{a}_i)$ is only 1 when $i = j$ and 0 otherwise. Then each \mathbf{a}_i is 1-independent to the predecessors on f_i and g_i , and thus the eluder dimension is lower bounded by $\binom{d}{p}$.
- \mathcal{F}_{EV} : the sequence is $\{\mathbf{a}_i = \frac{1}{\sqrt{2}}(\mathbf{e}_{i_1} + \mathbf{e}_{i_2}) : i = (i_1, i_2) \in [d]^2, i_1 \leq i_2\}$. For $f_j(\mathbf{a}) = \frac{1}{2} \mathbf{a}^\top (\mathbf{e}_{j_1} + \mathbf{e}_{j_2})(\mathbf{e}_{j_1} + \mathbf{e}_{j_2})^\top \mathbf{a}$ and $g_j(\mathbf{a}) = \frac{1}{2} \mathbf{a}^\top (\mathbf{e}_{j_1} - \mathbf{e}_{j_2})(\mathbf{e}_{j_1} - \mathbf{e}_{j_2})^\top \mathbf{a}$ with $j_1 \leq j_2$, $(f_j - g_j)(\mathbf{a}_i)$ is only 1 when $i = j$ and 0 otherwise. Then each \mathbf{a}_i is 1-independent to the predecessors on f_i and g_i , and thus the eluder dimension is lower bounded by $\binom{d}{2}$.
- \mathcal{F}_{LR} : the sequence is $\{\mathbf{A}_i = \frac{1}{2} \mathbf{e}_{i_1} \mathbf{e}_{i_2}^\top + \mathbf{e}_{i_1} \mathbf{e}_{i_2}^\top : i = (i_1, i_2) \in [d]^2, i_1 \leq i_2\}$. For $f_j(\mathbf{A}) = \langle \frac{1}{2}(\mathbf{e}_{j_1} \mathbf{e}_{j_2}^\top + \mathbf{e}_{j_2} \mathbf{e}_{j_1}^\top), \mathbf{A} \rangle$ with $j_1 \leq j_2$, $f_j(\mathbf{A}_i)$ is only 1 when $i = j$ and 0 otherwise. Then each \mathbf{A}_i is 1-independent to the predecessors on f_i and zero, and thus the eluder dimension is lower bounded by $\binom{d}{2}$.

□

Then we are all set for the results in the first line of 1. Notice that when we choose $\alpha = O(1/T^2)$ and $\epsilon = O(1/T^2)$ in our analysis of Algorithm 2, the regret upper bound would only expand by $\log(T)$ factors.

C Omitted Proofs for Quadratic Reward

In this section we include all the omitted proof of the theorems presented in the main paper.

C.1 Omitted Proofs of Main Results for Stochastic Bandit Eigenvector Problem

Proof of Theorem 3.3. Notice in Algorithm 1, for each iterate \mathbf{a} , its next iterate \mathbf{y} satisfies

$$\begin{aligned} \mathbf{y} &= \frac{1}{n_s} \sum_{i=1}^{n_s} (\mathbf{a}/2 + \mathbf{z}_i/2)^\top \mathbf{M}(\mathbf{a}/2 + \mathbf{z}_i/2) \mathbf{z}_i + \eta_i \mathbf{z}_i \\ &= \frac{m_s}{n_s} \sum_{i=1}^{n_s} \left(\frac{1}{4} \mathbf{a}^\top \mathbf{M} \mathbf{a} + \frac{1}{2} \mathbf{a}^\top \mathbf{M} \mathbf{z}_i + \eta_i \right) \mathbf{z}_i. \end{aligned}$$

Therefore $\mathbb{E}[\mathbf{y}] = \frac{1}{2}\mathbf{M}\mathbf{a}$. We can write $2\mathbf{y} = \mathbf{M}\mathbf{a} + \mathbf{g}$ where $\mathbf{g} := \frac{m_s}{n_s} \sum_{i=1}^{n_s} (\frac{1}{2}\mathbf{a}^\top \mathbf{M}\mathbf{a} + 2\eta_i)\mathbf{z}_i$.

With Claim D.12 and Claim D.11 we get that $\|\mathbf{g}\| \leq C\sqrt{\frac{m_s \log^2(n/\delta) \log(d/\delta)d}{n_s}}$. Therefore with our choice of $n_s \geq \tilde{\Theta}(\frac{d^2}{\varepsilon_s^2(\lambda_1 - |\lambda_2|)^2})$ we guarantee $\|\mathbf{g}\| \leq \varepsilon_s(\lambda_1 - |\lambda_2|)$. Therefore it satisfies the requirements for noisy power method, and by applying Corollary C.4, we have with $L = O(\kappa \log(d/\varepsilon))$ iterations we will be able to find $\|\hat{\mathbf{a}} - \mathbf{a}^*\| \leq \varepsilon$. By setting $\delta < 0.1/L$ in the algorithm we can guarantee the whole process succeed with high probability. Altogether it is sufficient to take $Ln_s = \tilde{O}(\kappa d^2/(\varepsilon\Delta)^2)$ actions to get an ε -optimal arm.

Finally to get the cumulative regret bound, we apply Claim D.7 with $A = \frac{d^2\kappa}{\Delta^2}$ and $a = 2$. Therefore we set $\varepsilon = A^{1/4}T^{-1/4} = \frac{d^{1/2}\kappa^{1/4}}{\Delta^{1/2}T^{1/4}}$ and get:

$$\text{Reg}(T) \lesssim T^{1/2}A^{1/2}r^* = \sqrt{\frac{d^2\kappa}{\Delta^2}}Tr^* = \sqrt{d^2\kappa^3T}.$$

□

Corollary C.1 (Formal statement for Corollary 3.6). *In Algorithm 1, by setting $\alpha = 1 - \varepsilon^2/2$, one can get ε -optimal reward with a total of $\tilde{O}(d^2\lambda_1^2/\varepsilon^4)$ total samples to get \mathbf{a} such that $r^* - f(\mathbf{a}) \leq \varepsilon$. Therefore one can get an accumulative regret of $\tilde{O}(\lambda_1^{3/5}d^{2/5}T^{4/5})$.*

Proof of Lemma 3.6. In order to find an arm with $\lambda_1\varepsilon^2$ -optimal reward, one will want to recover an arm that is $\varepsilon/2$ -close (meaning to find an \mathbf{a} such that $\tan\theta(\mathbf{V}_l, \mathbf{a}) \leq \varepsilon/2$) to the top eigenspace $\text{span}(\mathbf{v}_1, \dots, \mathbf{v}_l)$, where l satisfies $\lambda_l \geq \lambda_1 - \tilde{\varepsilon}$ and $\lambda_{l+1} \leq \lambda_1 - \tilde{\varepsilon}$. Here we set $\tilde{\varepsilon} := \lambda_1\varepsilon^2/2$. We first show 1) this is sufficient to get an $\lambda_1\varepsilon$ -optimal reward, and next show 2) how to set parameter to achieve this.

To get 1), we write $\mathbf{V}_l = [\mathbf{v}_1, \dots, \mathbf{v}_l] \in \mathbb{R}^{d \times l}$ and $\mathbf{V}_l^\perp = [\mathbf{v}_{l+1}, \dots, \mathbf{v}_k]$. When $\tan\theta(\mathbf{V}_l, \mathbf{a}_T) = \|\mathbf{V}_l^\perp \mathbf{a}\|/\|\mathbf{V}_l \mathbf{a}\| \leq \varepsilon/2$, from the proof of Claim D.6, we get $r^* - f(\mathbf{a}) \leq \min\{\lambda_1, \lambda_1 2(\varepsilon/2)^2 + \tilde{\varepsilon}\} = \lambda_1\varepsilon^2$.

Now to get 2), we note that in each iteration we try to conduct the power iteration to find an action $\tan\theta(\mathbf{V}_l, \hat{\mathbf{a}}) \leq \varepsilon/2$ and with eigengap $\geq \tilde{\varepsilon} := \lambda_1\varepsilon^2/2$. Therefore it is sufficient to let $\|\mathbf{g}\| \leq 0.1\tilde{\varepsilon}$ and $|\mathbf{v}_1^\top \mathbf{g}| \leq 0.1\tilde{\varepsilon}\frac{1}{\sqrt{d}}$, and thus $n_s \geq \tilde{\Theta}(\frac{d^2}{\varepsilon^2\tilde{\varepsilon}^2}) \leq \tilde{\Theta}(d^2/\lambda_1^2\varepsilon^6)$. Together we need $\lambda_1/\tilde{\varepsilon} \log(2d/\varepsilon)n_s = \tilde{\Theta}(d^2/\lambda_1^2\varepsilon^8)$ samples to get an $\lambda_1\varepsilon^2$ -optimal reward. Namely we get $\tilde{\varepsilon}$ -optimal reward with $\tilde{O}(d^2\lambda_1^2/\varepsilon^4)$ samples.

Finally by applying Claim D.7 we get:

$$\mathfrak{R}(T) \lesssim (d^2\lambda_1^2)^{\frac{1}{5}}T^{\frac{4}{5}}\lambda_1^{\frac{1}{5}} \leq \tilde{O}(\lambda_1^{3/5}d^{2/5}T^{4/5}).$$

□

Theorem C.2 (Formal statement of Theorem 3.7). *In Algorithm 3, if we set $n = \tilde{\Theta}(\frac{d^2\lambda_1^2}{\varepsilon^2\lambda_k^2})$, $m = d \log(n/\delta)$, $L = \Theta(\log(d/\varepsilon))$, $\delta = 0.1/L$, we will be able to identify an action $\hat{\mathbf{a}}$ that yield at most ε -regret with probability 0.9. Therefore by applying the standard PAC to regret conversion as discussed in Claim D.7 we get a cumulative regret of $\tilde{O}(\lambda_1^{1/3}k^{1/3}(\tilde{\kappa}dT)^{2/3})$ for large enough T , where $\tilde{\kappa} = \lambda_1/|\lambda_k|$.*

On the other hand, we set $n = \tilde{\Theta}(\frac{d^2k^2}{\varepsilon^2})$ and keep the other parameters. If we play Algorithm 3 k times by setting $k' = 2, 4, 6, \dots, 2k$ and select the best output among them, we can get a gap-free cumulative regret of $\tilde{O}(\lambda_1^{1/3}k^{4/3}(dT)^{2/3})$ for large enough T with high probability.

Proof of Theorem 3.7. First we show the first setting identify an ε -optimal reward with $\tilde{O}(\tilde{\kappa}^2d^2k\varepsilon^{-2})$ samples.

Similarly as Theorem 3.8, when setting $n \geq \tilde{\Theta}(d^2/(\sigma_k^2\tilde{\varepsilon}^2))$, we can find \mathbf{X}_L that satisfies $\|(\mathbf{X}_L\mathbf{X}_L^\top - I)\mathbf{U}\| \leq \tilde{\varepsilon}$, and therefore we recover an $\mathbf{Y}_L = \mathbf{M}\mathbf{X}_{L-1} + \mathbf{G}_L$ with $\|\mathbf{G}_L\| \leq \sigma_k\tilde{\varepsilon}$ and

Algorithm 3 Gap-free Subspace Iteration for Bilinear Bandit

- 1: **Input:** Quadratic reward $f : \mathcal{X} \rightarrow \mathbb{R}$ generating noisy reward, failure probability δ , error ϵ .
 - 2: **Initialization:** Set $k' = 2k$. Initial candidate matrix $\mathbf{X}_0 \in \mathbb{R}^{d \times k'}$, $\mathbf{X}_0(j) \in \mathbb{R}^d$, $j = 1, 2, \dots, k'$ is the j -th column of \mathbf{X}_0 and are i.i.d sampled on the unit sphere \mathbb{S}^{d-1} uniformly. Sample variance m , # sample per iteration n , total iteration L .
 - 3: **for** Iteration l from 1 to L **do**
 - 4: **for** s from 1 to k' **do**
 - 5: **Noisy subspace iteration:**
 - 6: Sample $\mathbf{z}_i \sim \mathcal{N}(0, 1/mI_d)$, $i = 1, 2, \dots, n_s$.
 - 7: Calculate tentative rank-1 arms $\tilde{\mathbf{a}}_i = \frac{1}{2}(\mathbf{X}_{l-1}(s) + \mathbf{z}_i)$.
 - 8: Conduct estimation $\mathbf{Y}_l(s) \leftarrow 4m/n \sum_{i=1}^n (f(\tilde{\mathbf{a}}_i) + \eta_i) \mathbf{z}_i$. ($\mathbf{Y}_l \in \mathbb{R}^{d \times k'}$)
 - 9: Let $\mathbf{Y}_l = \mathbf{X}_l \mathbf{R}_l$ be a QR-factorization of \mathbf{Y}_l
 - 10: Update target arm $\mathbf{a}_l \leftarrow \arg \max_{\|\mathbf{a}\|=1} \mathbf{a}^\top \mathbf{Y}_l \mathbf{X}_{l-1}^\top \mathbf{a}$.
 - 11: **Output:** \mathbf{a}_L .
-

$\|\mathbf{Y}_L \mathbf{X}_{L-1}^\top - \mathbf{M}\|_2 = \|\mathbf{M} \mathbf{X}_{L-1} \mathbf{X}_{L-1}^\top - \mathbf{M} + \mathbf{G}_L \mathbf{X}_{L-1}^\top\|_2 \leq (\lambda_1 + |\lambda_k|) \tilde{\epsilon}$. Therefore by definition of \mathbf{a}_L , $\mathbf{a}_L^\top \mathbf{Y}_L \mathbf{X}_{L-1}^\top \mathbf{a}_L = \max_{\|\mathbf{a}\|=1} \mathbf{a}^\top (\mathbf{M} \mathbf{X}_{L-1} \mathbf{X}_{L-1}^\top + \mathbf{G}_L \mathbf{X}_{L-1}^\top) \mathbf{a} \geq \lambda_1 - (\lambda_1 + |\lambda_k|) \tilde{\epsilon}$. Therefore $\mathbf{a}_L^\top \mathbf{M} \mathbf{a}_L \geq \lambda_1 - 2(\lambda_1 + |\lambda_k|) \tilde{\epsilon}$. Therefore we set $2(\lambda_1 + |\lambda_k|) \tilde{\epsilon} = \epsilon$, i.e., $\tilde{\epsilon} = 0.5\epsilon / (\lambda_1 + |\lambda_k|)$ which will get a total sample of $T = \tilde{\Theta}(kn) = \tilde{\Theta}(d^2 \tilde{\kappa}^2 k \epsilon^{-2})$. Then by applying Claim D.7 we get the cumulative regret bound.

Next we show how to estimate the action with $\tilde{O}(d^2 k^4 \epsilon^{-2})$ samples. To achieve this result, we need to slightly alter Algorithm 3 where we respectively set $k' = 2, 4, 6, \dots, 2k$ and keep the best arm among the k outputs. We argue that among all the choices of k' , at least for one $l \in [k]$, $k' = 2l$, we have $|\lambda_l| - |\lambda_{l+1}| \geq \lambda_l/k$. Notice with similar argument as above, when we set $n = \tilde{\Theta}(d^2 \lambda_l^{-2} \tilde{\epsilon}^{-2}) \leq \tilde{\Theta}(d^2 k^2 \lambda_1^{-2} \tilde{\epsilon}^{-2})$ we can get $\|\mathbf{G}\| \leq \tilde{\epsilon} \lambda_l$ as required by Corollary C.4, the total number of iterations $L = O(\sigma_l / (\sigma_l - \sigma_{l+1}) \log(2d/\epsilon)) = \tilde{O}(k)$. Finally by setting $\tilde{\epsilon} = \epsilon / (4\lambda_1)$ we get the overall samples we required is $\tilde{O}(k^2 n) = \tilde{O}(d^2 k^4 \epsilon^{-2})$.

For both settings, directly applying our arguments in the PAC to regret conversion: Claim D.7 will finish the proof. \square

C.2 Omitted Details of Main Results of Low-Rank Linear Reward

Algorithm 4 Subspace Iteration Exploration for Low-rank Linear Reward.

- 1: **Input:** Quadratic function $f : \mathcal{A} \rightarrow \mathbb{R}$ with noisy reward, failure probability δ , error ϵ .
 - 2: **Initialization:** Set $k' = 2k$. Initial candidate matrix $\mathbf{X}_0 \in \mathbb{R}^{d \times k'}$, $\mathbf{X}_0(j) \in \mathbb{R}^d$, $j = 1, 2, \dots, k'$ is the j -th column of \mathbf{X}_0 and are i.i.d sampled on the unit sphere \mathbb{S}^{d-1} uniformly. Sample variance m , # sample per iteration n , total iteration L .
 - 3: **for** Iteration l from 1 to L **do**
 - 4: Sample $\mathbf{z}_i \sim \mathcal{N}(0, 1/mI_d)$, $i = 1, 2, \dots, n$.
 - 5: **for** s from 1 to k' **do**
 - 6: **Noisy subspace iteration:**
 - 7: Calculate tentative rank-1 actions $\tilde{\mathbf{A}}_i = \mathbf{X}_{l-1}(s) \mathbf{z}_i^\top$.
 - 8: Conduct estimation $\mathbf{Y}_l(s) \leftarrow m/n \sum_{i=1}^n (\langle \mathbf{M}, \tilde{\mathbf{A}}_i \rangle + \eta_{i,s}) \mathbf{z}_i$. ($\mathbf{Y}_l \in \mathbb{R}^{d \times k'}$)
 - 9: Let $\mathbf{Y}_l = \mathbf{X}_l \mathbf{R}_l$ be a QR-factorization of \mathbf{Y}_l
 - 10: Update target action $\mathbf{A}_l \leftarrow \mathbf{Y}_l \mathbf{X}_l^\top$.
 - 11: **Output:** $\hat{\mathbf{A}} = \mathbf{A}_L / \|\mathbf{A}_L\|_F$
-

Theorem C.3 (Formal statement of Theorem 3.8). *In Algorithm 4, for large enough constants C_n, C_L, C_m , let $n = C_n d^2 \log^2(d/\delta) \sigma_k^{-2} \epsilon^{-2}$, $m = C_m d \log(n/\delta)$, and $L = C_L \log(d/\epsilon)$, \mathbf{X}_L*

satisfies $\|(\mathbf{I} - \mathbf{X}_L \mathbf{X}_L^\top) \mathbf{V}\| \leq \varepsilon/4$, and the output $\widehat{\mathbf{A}}$ satisfies $\|\widehat{\mathbf{A}} - \mathbf{A}^*\|_F \leq \|\mathbf{M}\|_F \varepsilon$. Altogether to get an ε -optimal action, it is sufficient to have total sample complexity of $T \leq \widetilde{O}(d^2 k \lambda_k^{-2} \varepsilon^{-2})$.

Proof of Theorem 3.8. Let $\mathbf{M} = \mathbf{V} \Sigma \mathbf{V}^\top$. From Claim C.6 we get that for each noisy subspace iteration step we get $\mathbf{Y}_l = \mathbf{M} \mathbf{X}_l + \mathbf{G}_l$ with $5\|\mathbf{G}_l\| \leq \varepsilon \sigma_k$ and $\|\mathbf{V}^\top \mathbf{G}\| \leq \sigma_k \sqrt{k}/3\sqrt{d} \leq \sigma_k(\sqrt{2k} - \sqrt{k})/2\sqrt{d}$. Therefore we can apply Corollary C.4, and get $\|\mathbf{V}(\mathbf{X}_L \mathbf{X}_L^\top - \mathbf{I})\| \leq \varepsilon/4$ with $O(\log 2d/\varepsilon)$ steps. Therefore we have:

$$\begin{aligned} \|\mathbf{A}_L - \mathbf{M}\|_F &= \|(\mathbf{M} \mathbf{X}_L + \mathbf{G}_L) \mathbf{X}_L^\top - \mathbf{M}\|_F = \|\mathbf{V}^\top \Sigma \mathbf{V}(\mathbf{X}_L \mathbf{X}_L^\top - \mathbf{I}) + \mathbf{G}_L \mathbf{X}_L^\top\|_F \\ &\leq \|\mathbf{M}\|_F \|\mathbf{V}(\mathbf{X}_L \mathbf{X}_L^\top - \mathbf{I})\| + \|\mathbf{G}_L\| \|\mathbf{X}_L\|_F \\ &\leq (\|\mathbf{M}\|_F + \sigma_k) \varepsilon/4 < \|\mathbf{M}\|_F \varepsilon/2. \end{aligned}$$

Meanwhile, notice $\|\mathbf{A}^*\|_F = 1$, $\|\mathbf{M}\|_F = r^*$ and $\|\widehat{\mathbf{A}}\|_F = 1$. $\|\mathbf{A}_L/r^* - \mathbf{A}^*\|_F \leq \varepsilon/2$. $\|\widehat{\mathbf{A}} - \mathbf{A}^*\|_F = \|\mathbf{A}_L/\|\mathbf{A}_L\|_F - \mathbf{A}^*\|_F = \|\text{vec}(\mathbf{A}_L)/\|\text{vec}(\mathbf{A}_L)\|_2 - \text{vec}(\mathbf{A}^*)\|_2$.

Write $\theta_A := \theta(\text{vec}(\mathbf{A}_L), \text{vec}(\mathbf{A}^*))$. The worst case that makes $\|\text{vec}(\widehat{\mathbf{A}}) - \text{vec}(\mathbf{A}^*)\|$ to be larger than $\|\text{vec}(\mathbf{A}_L/r^*) - \text{vec}(\mathbf{A}^*)\|$ is when $\|\text{vec}(\mathbf{A}_L/r^*) - \text{vec}(\mathbf{A}^*)\| = \sin \theta_A$ and $\|\text{vec}(\widehat{\mathbf{A}}) - \text{vec}(\mathbf{A}^*)\|$ is always $2 \sin(\theta_A/2)$. Notice trivially $2 \sin(\theta_A/2) \leq 2 \sin(\theta_A)$ Therefore we could get $\|\widehat{\mathbf{A}} - \mathbf{A}^*\|_F \leq 2\|\mathbf{A}_L/r^* - \mathbf{A}^*\|_F \leq \varepsilon$. □

Proof of Corollary 3.9. The corollary uses a special property of the strongly convex action set that ensures: $\mathbf{A}^* = \mathbf{M}/r^*$. With $\widehat{\mathbf{A}}$ that satisfies $\|\widehat{\mathbf{A}}\|_F = 1$, we have

$$\begin{aligned} r^* - f_{\mathbf{M}}(\mathbf{A}) &= r^* - \langle \widehat{\mathbf{A}}, \mathbf{M} \rangle = r^* - \langle \widehat{\mathbf{A}}, r^* \mathbf{A}^* \rangle \\ &= \frac{r^*}{2} (2 - 2\langle \widehat{\mathbf{A}}, \mathbf{A}^* \rangle) = \frac{r^*}{2} (\|\widehat{\mathbf{A}}\|_F^2 + \|\mathbf{A}^*\|_F^2 - 2\langle \widehat{\mathbf{A}}, \mathbf{A}^* \rangle) \\ &= \frac{r^*}{2} \|\widehat{\mathbf{A}} - \mathbf{A}^*\|_F^2 \leq \frac{r^* \varepsilon^2}{2} \end{aligned} \tag{2}$$

Therefore, with first $T_1 = \widetilde{O}(d^2 k \lambda_k^{-2} \varepsilon^{-2})$ exploratory samples we get $r^* - f(\widehat{\mathbf{A}}) \leq r^* \varepsilon^2/2 = r^* \sqrt{\frac{d^2 k}{\lambda_k^2 T}} = \sqrt{\frac{(r^*)^2 d^2 k}{\lambda_k^2 T}}$. Together we have:

$$\begin{aligned} \mathfrak{R}(T) &= \sum_{t=1}^{T_1} r^* - f(\mathbf{A}_t) + \sum_{t=T_1+1}^T r^* - f(\widehat{\mathbf{A}}) \\ &< r^* T_1 + T r^* \varepsilon^2 \\ &\leq \widetilde{O}(\sqrt{d^2 k (r^*)^2 \lambda_k^{-2} T}). \end{aligned}$$

□

Proof of Theorem 3.10. We find an l to be the smallest integer such that $\sum_{i=l+1}^k \sigma_i^2 \leq \varepsilon^2 \|\mathbf{M}\|_F^2$. Then we have $\sigma_l \geq \varepsilon/\sqrt{k-l} > \varepsilon/\sqrt{k}$.

Notice that in Algorithm 4, we set $n \geq \widetilde{\Theta}(\frac{d^2 k}{(r^*)^2 \varepsilon^4})$ large enough such that $\|\mathbf{G}\|_2 \leq O(\|\mathbf{M}\|_F \varepsilon^2 / \sqrt{k}) \lesssim \varepsilon(\sigma_l - 0)$ and $\|\mathbf{U}^\top \mathbf{G}\|_2 \leq \|\mathbf{M}\|_F \varepsilon / \sqrt{k} \frac{\sqrt{k^l - \sqrt{k-1}}}{2\sqrt{d}}$. (This comes from the argument proved in Claim C.6.)

Therefore by conducting noisy power method we get with $O(nk) = \widetilde{O}(\frac{d^2 k^2}{(r^*)^2 \varepsilon^4})$ samples we can get an action $\widehat{\mathbf{A}}$ that satisfies:

$$\|\mathbf{M} - \mathbf{X}_L \mathbf{X}_L^\top \mathbf{M}\|_F^2 \leq \sum_{i=l+1}^k \sigma_i^2 + l \varepsilon^2 \sigma_l^2 \leq 2\|\mathbf{M}\|_F^2 \varepsilon^2.$$

Therefore we could get $\|\mathbf{A}^* - \widehat{\mathbf{A}}\| \leq 2\epsilon$, and with similar argument as (2) we have $r^* - f(\widehat{\mathbf{A}}) \leq \|\mathbf{M}\|_F \epsilon^2$.

Therefore if we want to take a total of T actions, we will set $\epsilon^6 = \widetilde{\Theta}(\frac{d^2 k^2}{(r^*)^{2T}})$ and we get:

$$\begin{aligned} \mathfrak{R}(T) &= \sum_{t=1}^{T_1} r^* - f(\mathbf{A}_t) + \sum_{t=T_1+1}^T r^* - f(\widehat{\mathbf{A}}) \\ &< r^* T_1 + T r^* \epsilon^2 \\ &\leq \widetilde{O}(d^{2/3} k^{2/3} (r^*)^{1/3} T^{2/3}). \end{aligned}$$

□

C.3 Technical Details for Quadratic Reward

Noisy Power Method.

Corollary C.4 (Adapted from Corollary 1.1 from [37]). *Let $k' \geq l$. Let $\mathbf{U} \in \mathbb{R}^{d \times l}$ represent the top l singular vectors of \mathbf{M} and let $\sigma_1 \geq \dots \geq \sigma_k > 0$ denote its singular values. Suppose \mathbf{X}_0 is an orthonormal basis of a random k' -dimensional subspace. Further suppose that at every step of NPM we have*

$$\begin{aligned} 5\|\mathbf{G}\| &\leq \epsilon(\sigma_l - \sigma_{l+1}), \\ \text{and } 5\|\mathbf{U}^\top \mathbf{G}\| &\leq (\sigma_l - \sigma_{l+1}) \frac{\sqrt{k'} - \sqrt{l-1}}{2\sqrt{d}} \end{aligned}$$

for some fixed parameter $\epsilon < 1/2$. Then with all but $2^{-\Omega(k'+1-l)} + e^{\Omega(d)}$ probability, there exists an $L = O(\frac{\sigma_l}{\sigma_l - \sigma_{l+1}} \log(2d/\epsilon))$ so that after L steps we have that $\|(\mathbf{I} - \mathbf{X}_L \mathbf{X}_L^\top) \mathbf{U}\| \leq \epsilon$.

Theorem C.5 (Adapted from Theorem 2.2 from [11]). *Let $\mathbf{U}_l \in \mathbb{R}^{d \times l}$ represent the top l singular vectors of \mathbf{M} and let $\sigma_1 \geq \dots \geq \sigma_k > 0$ denote its singular values. Naturally $l \leq k$. Suppose \mathbf{X}_0 is an orthonormal basis of a random k' -dimensional subspace where $k' \geq k$. Further suppose that at every step of NPM we have*

$$\begin{aligned} \|\mathbf{G}\| &\leq O(\epsilon \sigma_l), \\ \text{and } \|\mathbf{U}_k^\top \mathbf{G}\|_2 &\leq O(\sigma_l \frac{\sqrt{k'} - \sqrt{k-1}}{2\sqrt{d}}) \end{aligned}$$

for small enough ϵ . Then with all but $2^{-\Omega(k'+1-k)} + e^{\Omega(d)}$ probability, there exists an $L = O(\log(2d/\epsilon))$ so that after L steps we have that $\|(\mathbf{I} - \mathbf{X}_L \mathbf{X}_L^\top) \mathbf{U}_l\| \leq \epsilon$. Furthermore:

$$\|\mathbf{M} - \mathbf{X}_L \mathbf{X}_L^\top \mathbf{M}\|_F^2 \leq \sum_{i=l+1}^k \sigma_i^2 + l \sigma_l^2.$$

Concentration Bounds.

Claim C.6. *Write the eigendecomposition for \mathbf{M} as $\mathbf{M} = \mathbf{U} \boldsymbol{\Sigma} \mathbf{U}^\top$. In Algorithm 4, when $n \geq \widetilde{\Theta}(d^2 / (\lambda_k^2 \epsilon^2))$, the noisy subspace iteration step can be written as: $\mathbf{Y}_l = \mathbf{M} \mathbf{X}_{l-1} + \mathbf{G}_l$, where the noise term satisfies:*

$$\begin{aligned} 5\|\mathbf{G}_l\| &\leq \epsilon |\lambda_k| \\ 5\|\mathbf{U}^\top \mathbf{G}_l\| &\leq \epsilon |\lambda_k| \frac{\sqrt{k}}{3\sqrt{d}}. \end{aligned}$$

with high probability for our choice of n .

Proof. For compact notation, write vector $\boldsymbol{\eta}_i := [\eta_{i,1}, \eta_{i,2}, \dots, \eta_{i,k'}]^\top \in \mathbb{R}^{k'}$. We have:

$$\begin{aligned} \mathbf{G}_l(s) &= \frac{m}{n} \sum_{i=1}^n (\mathbf{z}_i^\top \mathbf{M} \mathbf{X}_l(s)) \mathbf{z}_i + \frac{m}{n} \sum_{i=1}^n \eta_{i,s} \mathbf{z}_i - \mathbf{M} \mathbf{X}_l(s), \text{ therefore} \\ \mathbf{G}_l &= \left(\frac{m}{n} \sum_{i=1}^n [\mathbf{z}_i \mathbf{z}_i^\top] - \mathbf{I} \right) \mathbf{M} \mathbf{X}_l + \frac{m}{n} \sum_{i=1}^n \mathbf{z}_i \boldsymbol{\eta}_i^\top. \end{aligned}$$

First note that for orthogonal matrix \mathbf{X}_l , $\|\mathbf{M}\mathbf{X}_l\| \leq \lambda_1$, and $\|\frac{m}{n} \sum_{i=1}^n [\mathbf{z}_i \mathbf{z}_i^\top] - I\| \leq O(\sqrt{\frac{d+\log(1/\delta)}{n}})$. The bottleneck is from the second term and we will use Matrix Bernstein to concentrate it. Write $\mathbf{S}_i = \frac{m}{n} \mathbf{z}_i \boldsymbol{\eta}_i^\top$. We have $\|\mathbf{S}_i\| \leq O(\frac{\sqrt{mk'} \log(n/\delta)}{n})$ with probability $1 - \delta$ and $\mathbb{E}[\sum_i \mathbf{S}_i \mathbf{S}_i^\top] = \frac{mk'}{n} I_d$ and $\mathbb{E}[\sum_i \mathbf{S}_i^\top \mathbf{S}_i] = \frac{md}{n} I_{k'}$. Therefore with matrix Bernstein we can get that $\|\sum_i \mathbf{S}_i\|_i \leq O(\sqrt{\frac{md}{n}} \log(d/\delta))$ with probability $1 - \delta$.

Therefore for $n \geq \tilde{\Omega}(d^2/(\lambda_k^2 \varepsilon^2))$, we can get that $5\|\mathbf{G}_l\| \leq \varepsilon|\lambda_k|$.

Similarly since $\mathbf{U}^\top \mathbf{z}_i \sim \mathcal{N}(0, \frac{1}{m} I_{k'})$, with the same argument one can easily get that $\|\mathbf{U}^\top \mathbf{G}_l\| \leq O(\sqrt{\frac{mk'}{n}} \log(d/\delta))$. Therefore with the same lower bound for n one can get $15\|\mathbf{U}^\top \mathbf{G}_l\| \leq \varepsilon|\lambda_k| \sqrt{\frac{k}{d}}$. \square

C.4 Omitted Proof for RL with Quadratic Q function

Algorithm 5 Learn policy complete polynomial with simulator.

- 1: **Initialize:** Set $n = \tilde{\Theta}(\tilde{\kappa}^2 d^2 H^3 / \varepsilon^2)$, Oracle to estimate \hat{T}_h from noisy observations.
 - 2: **for** $h = H, \dots, 1$ **do**
 - 3: Sample $\phi(s_h^i, a_h^i), i \in [n]$ from standard Gaussian $N(0, I_d)$
 - 4: **for** $i \in [n]$ **do**
 - 5: Query (s_h^i, a_h^i) and use π_{h+1}, \dots, π_H as the roll-out to get estimation $\hat{Q}_h^{\pi_{h+1}, \dots, \pi_H}(s_h^i, a_h^i)$
 - 6: Retrieve $\hat{\mathbf{M}}_h$ from estimation $\hat{Q}_h^{\pi_{h+1}, \dots, \pi_H}(s_h^i, a_h^i), i \in [n]$
 - 7: Set $\hat{Q}_h(s, a) \leftarrow f_{\hat{T}_h}$
 - 8: Set $\pi_h(s) \leftarrow \arg \max_{a \in \mathcal{S}} \hat{Q}_h(s, a)$
 - 9: **Return** π_1, \dots, π_H
-

Proof of Theorem 3.13. With the oracle, at horizon H , we can estimate $\hat{\mathbf{M}}_H$ that is ε/H close to \mathbf{M}_H^* in spectral norm through noisy observations from the reward function with $\tilde{O}(\tilde{\kappa}^2 d^2 H^2 / \varepsilon^2)$ samples. Next, for each horizon $h = H - 1, H - 1, \dots, 1$, sample $s'_i \sim \mathbb{P}(\cdot | s, a)$, we define $\eta_i = \max_{a'} f_{\hat{\mathbf{M}}_{h+1}}(s'_i, a') - \mathbb{E}_{s' \sim \mathbb{P}(\cdot | s, a)} \max_{a'} f_{\hat{\mathbf{M}}_{h+1}}(s', a')$. η_i is mean-zero and $O(1)$ -sub-gaussian since it is bounded. Denote \mathbf{M}_h as the matrix that satisfies $f_{\mathbf{M}_h} := \mathcal{T} f_{\hat{\mathbf{M}}_{h+1}}$, which is well-defined due to Bellman completeness. We estimate $\hat{\mathbf{M}}_h$ from the noisy observations $y_i = r_h(s, a) + \max_{a'} f_{\hat{\mathbf{M}}_{h+1}}(s'_i, a') = \mathcal{T} f_{\hat{\mathbf{M}}_{h+1}} + \eta_i =: f_{\mathbf{M}_h} + \eta_i$. Therefore with the oracle, we can estimate $\hat{\mathbf{M}}_h$ such that $\|\hat{\mathbf{M}}_h - \mathbf{M}_h\|_2 \leq \varepsilon/H$ with $\Theta(\tilde{\kappa}^2 d^2 k^2 H^2 / \varepsilon^2)$ bandits. Together we have:

$$\begin{aligned}
\|f_{\hat{\mathbf{M}}_h} - f_{\mathbf{M}_h^*}\|_\infty &= \|\hat{\mathbf{M}}_h - \mathbf{M}_h^*\| \\
&\leq \|\hat{\mathbf{M}}_h - \mathbf{M}_h\| + \|\mathbf{M}_h - \mathbf{M}_h^*\| \\
&\leq \varepsilon/H + \|\mathcal{T} f_{\hat{\mathbf{M}}_{h+1}} - \mathcal{T} f_{\mathbf{M}_{h+1}^*}\|_\infty \\
&\leq \varepsilon/H + \|f_{\hat{\mathbf{M}}_{h+1}} - f_{\mathbf{M}_{h+1}^*}\|_\infty \\
&\leq 2\varepsilon/H + \|f_{\hat{\mathbf{M}}_{h+2}} - f_{\mathbf{M}_{h+2}^*}\|_\infty \\
&\leq \dots \\
&\leq (H - h)\varepsilon/H.
\end{aligned}$$

Finally for $h = 1$ we have $\|\hat{\mathbf{M}}_1 - \mathbf{M}^*\| \leq \varepsilon$ if we sample $n = \tilde{\Theta}(\tilde{\kappa}^2 d^2 k^2 H^2 / \varepsilon^2)$ for each $h \in [H]$. Therefore for all the H timesteps we need $\Theta(\tilde{\kappa}^2 d^2 k^2 H^3 / \varepsilon^2)$. \square

D Technical details for General Tensor Reward

Algorithm 6 Phased elimination with zeroth order exploration.

- 1: **Input:** Function $f : \mathcal{A} \rightarrow \mathbb{R}$ of polynomial degree p generating noisy reward, failure probability δ , error ε .
 - 2: **Initialization:** $L_0 = C_L k \log(1/\delta)$; Total number of stages $S = C_S \lceil \log(1/\varepsilon) \rceil + 1$, $\mathcal{A}_0 = \{\mathbf{a}_0^{(1)}, \mathbf{a}_0^{(2)}, \dots, \mathbf{a}_0^{(L_0)}\}$ where each $\mathbf{a}_0^{(l)}$ is uniformly sampled on the unit sphere \mathbb{S}^{d-1} . $\tilde{\varepsilon}_0 = 1$.
 - 3: **for** s from 1 to S **do**
 - 4: $\tilde{\varepsilon}_s \leftarrow \tilde{\varepsilon}_{s-1}/2$, $n_s \leftarrow C_n d^p \log(d/\delta) / \lambda_1^2 \tilde{\varepsilon}_s^2$, $n_s \leftarrow n_s \cdot \log^3(n_s/\delta)$, $m_s \leftarrow C_m d \log(n_s/\delta)$, $\mathcal{A}_s = \emptyset$.
 - 5: **for** l from 1 to L_{s-1} **do**
 - 6: **Zeroth-order optimization:**
 - 7: Locate current action $\tilde{\mathbf{a}} = \mathbf{a}_{s-1}^{(l)}$.
 - 8: **for** $\lceil (1/(1-\alpha)) \log(2d) \rceil$ times **do**
 - 9: Sample $\mathbf{z}_i \sim \mathcal{N}(0, 1/m_s I_d)$, $i = 1, 2, \dots, n_s$.
 - 10: Take actions $\mathbf{a}_i = (1 - \frac{1}{2p})\tilde{\mathbf{a}} + \frac{1}{2p}\mathbf{z}_i$ and observe $r_i = \mathbf{T}(\mathbf{a}_i) + \eta_i$, $i \in [n_s]$; Take actions $\frac{1}{2p}\mathbf{z}_i$ and observe $r'_i = \mathbf{T}(\frac{1}{2p}\mathbf{z}_i) + \eta'_i$, $i \in [n_s]$.
 - 11: Conduct estimation $\mathbf{y} \leftarrow 1/n_s \sum_{i=1}^{n_s} (r_i - r'_i)\mathbf{z}_i$.
 - 12: Update the current action $\tilde{\mathbf{a}} \leftarrow \mathbf{y} / \|\mathbf{y}\|$.
 - 13: Estimate the expected reward for $\tilde{\mathbf{a}}$ through n_s samples: $r_n(\tilde{\mathbf{a}}) = 1/n_s \sum_{i=1}^{n_s} (\mathbf{T}(\tilde{\mathbf{a}}) + \eta_i)$.
 - 14: **Candidate Elimination:**
 - 15: **if** $r_n \geq \lambda_1(1 - p\tilde{\varepsilon}_s^2)$ **then**
 - 16: Keep the action $\mathcal{A}_s \leftarrow \mathcal{A}_s \cup \{\tilde{\mathbf{a}}\}$
 - 17: Label the actions: $L_s = |\mathcal{A}_s|$, $\mathcal{A}_s =: \{\mathbf{a}_s^{(1)}, \dots, \mathbf{a}_s^{(L_s)}\}$.
 - 18: Run UCB (Algorithm 7) with the candidate set \mathcal{A}_S .
-

D.1 Technical Details for Symmetric Setting

Lemma D.1 (Zeroth order optimization for noiseless setting). *For $p \geq 3$, suppose $0.5\mathbf{a}^\top \mathbf{v}_1 > |\mathbf{a}^\top \mathbf{v}_j|$ for all $j \geq 2$, we have:*

$$\tan \theta(G(\mathbf{a}), \mathbf{v}_1) \leq \frac{1}{2} \tan \theta(\mathbf{a}, \mathbf{v}_1).$$

Proof. We first simplify $G(\mathbf{a}) = \sum_{j=1}^r \lambda_j \mathbf{v}_j \cdot S_j$, where

$$\begin{aligned} G(\mathbf{a}) &= \sum_{s=0}^{\lfloor (p-3)/2 \rfloor} \frac{(1 - \frac{1}{2p})^{p-2s-1} (\frac{1}{2p})^{2s+1}}{m^s} \binom{p}{2s+1} T(I^{\otimes s+1} \otimes \mathbf{a}^{\otimes p-2s-1}) \\ &= \sum_{s=0}^{\lfloor (p-3)/2 \rfloor} \frac{(1 - \frac{1}{2p})^{p-2s-1} (\frac{1}{2p})^{2s+1}}{m^s} \binom{p}{2s+1} \sum_{j=1}^k \lambda_j (\mathbf{v}_j^\top \mathbf{a})^{p-2j-1} \mathbf{v}_j \\ &= \sum_{j=1}^k \mathbf{v}_j \cdot \lambda_j \underbrace{\sum_{s=0}^{\lfloor (p-3)/2 \rfloor} \frac{(1 - \frac{1}{2p})^{p-2s-1} (\frac{1}{2p})^{2s+1}}{m^s} \binom{p}{2s+1}}_{S_j :=} (\mathbf{v}_j^\top \mathbf{a})^{p-2s-1} \\ &= \sum_{j=1}^k S_j \mathbf{v}_j. \end{aligned}$$

Notice for even p ,

$$\begin{aligned}
S_j &= \lambda_j (\mathbf{v}_j^\top \mathbf{a})^3 \cdot \sum_{s=0}^{p/2-2} \frac{(1 - \frac{1}{2p})^{p-2s-1} (\frac{1}{2p})^{2s+1}}{m^s} \binom{p}{2s+1} (\mathbf{v}_j^\top \mathbf{a})^{p-2s-4} \\
&= \lambda_j (\mathbf{v}_j^\top \mathbf{a})^3 \cdot \sum_{r=0}^{p/2-2} \frac{(1 - \frac{1}{2p})^{2r+3} (\frac{1}{2p})^{p-3-2r}}{m^{p/2-2-r}} \binom{p}{p-2r-3} (\mathbf{v}_j^\top \mathbf{a})^{2r}. \\
&\hspace{15em} (\text{let } 2r = p - 4 - 2s) \\
\frac{S_j}{\lambda_j (\mathbf{v}_j^\top \mathbf{a})^3} &= \sum_{r=0}^{p/2-2} \frac{(1 - \frac{1}{2p})^{2r+3} (\frac{1}{2p})^{p-3-2r}}{m^{p/2-2-r}} \binom{p}{p-2r-3} (\mathbf{v}_j^\top \mathbf{a})^{2r} \\
&\hspace{15em} (\text{Divide both sides by } \lambda_j (\mathbf{v}_j^\top \mathbf{a})^3) \\
&\leq \sum_{r=0}^{p/2-2} \frac{(1 - \frac{1}{2p})^{2r+3} (\frac{1}{2p})^{p-3-2r}}{m^{p/2-2-r}} \binom{p}{p-2r-3} (\mathbf{v}_1^\top \mathbf{a})^{2r}. \\
&\hspace{15em} (\text{Since the first term is constant and } |\mathbf{v}_j^\top \mathbf{a}| \leq \mathbf{v}_1^\top \mathbf{a} \text{ for } r \geq 1) \\
&= \frac{S_1}{\lambda_1 (\mathbf{v}_1^\top \mathbf{a})^3}.
\end{aligned}$$

Therefore for even $p \geq 4$:

$$|S_j| \leq \frac{|\lambda_j| |\mathbf{v}_j^\top \mathbf{a}|^3}{\lambda_1 |\mathbf{v}_1^\top \mathbf{a}|^3} S_1 \leq \frac{1}{4} \frac{|\mathbf{v}_j^\top \mathbf{a}|}{|\mathbf{v}_1^\top \mathbf{a}|} S_1, \forall j \geq 2. \quad (3)$$

Similarly for odd p , we have:

$$\begin{aligned}
S_j &= \lambda_j (\mathbf{v}_j^\top \mathbf{a})^2 \cdot \sum_{s=0}^{(p-3)/2} \frac{(1 - \frac{1}{2p})^{p-2s-1} (\frac{1}{2p})^{2s+1}}{m^s} \binom{p}{2s+1} (\mathbf{v}_j^\top \mathbf{a})^{p-2s-3} \\
&= \lambda_j (\mathbf{v}_j^\top \mathbf{a})^2 \cdot \sum_{r=0}^{(p-3)/2} \frac{(1 - \frac{1}{2p})^{2r+2} (\frac{1}{2p})^{p-2-2r}}{m^{(p-3)/2-r}} \binom{p}{p-2-2r} (\mathbf{v}_j^\top \mathbf{a})^{2r}, \\
&\hspace{15em} (\text{Let } r = (p-3)/2 - s) \\
\frac{S_j}{\lambda_j (\mathbf{v}_j^\top \mathbf{a})^2} &= \sum_{r=0}^{(p-3)/2} \frac{(1 - \frac{1}{2p})^{2r+2} (\frac{1}{2p})^{p-2-2r}}{m^{(p-3)/2-r}} \binom{p}{p-2-2r} (\mathbf{v}_j^\top \mathbf{a})^{2r} \\
&\hspace{15em} (\text{Divide both sides by } \lambda_j (\mathbf{v}_j^\top \mathbf{a})^2) \\
&\leq \sum_{r=0}^{(p-3)/2} \frac{(1 - \frac{1}{2p})^{2r+2} (\frac{1}{2p})^{p-2-2r}}{m^{(p-3)/2-r}} \binom{p}{p-2-2r} (\mathbf{v}_1^\top \mathbf{a})^{2r} \\
&\hspace{15em} (\text{Since the first term is constant and } |\mathbf{v}_j^\top \mathbf{a}| \leq \mathbf{v}_1^\top \mathbf{a} \text{ for } r \geq 1) \\
&= \frac{S_1}{\lambda_1 (\mathbf{v}_1^\top \mathbf{a})^2}.
\end{aligned}$$

Therefore for odd p we have:

$$|S_j| \leq \frac{|\lambda_j| |\mathbf{v}_j^\top \mathbf{a}|^2}{\lambda_1 |\mathbf{v}_1^\top \mathbf{a}|^2} S_1 \leq \frac{1}{2} \frac{|\mathbf{v}_j^\top \mathbf{a}|}{|\mathbf{v}_1^\top \mathbf{a}|} S_1, \forall j \geq 2. \quad (4)$$

Write $\mathbf{V} = [\mathbf{v}_2, \mathbf{v}_3, \dots, \mathbf{v}_k] \in \mathbb{R}^{d \times k}$ be the complement for \mathbf{v}_1 . Therefore for any \mathbf{x} without normalization, one can conveniently represent $|\tan \theta(\mathbf{x}, \mathbf{v}_1)|$ as $\|\mathbf{V}^\top \mathbf{x}\|_2 / |\mathbf{v}_1^\top \mathbf{x}|$.

$$\|\mathbf{V}^\top G(\mathbf{a})\|^2 = \sum_{j=2}^k S_j^2 \quad (5)$$

$$\leq \sum_{j=2}^k \frac{|\mathbf{v}_j^\top \mathbf{a}|^2}{4|\mathbf{v}_1^\top \mathbf{a}|^2} S_1^2 \quad (\text{from (4),(3)})$$

$$= \frac{1}{4} \tan^2 \theta(\mathbf{v}_1, \mathbf{a}) (\mathbf{v}_1^\top G(\mathbf{a}))^2. \quad (6)$$

Therefore for $p \geq 3$, $\tan \theta(G(\mathbf{a}), \mathbf{v}_1) \leq \frac{1}{2} \tan \theta(\mathbf{a}, \mathbf{v}_1)$. \square

D.1.1 Proof Sketch of Theorem 3.14

Definition D.2 (Zeroth order gradient function). For some scalar m , we define an empirical operator $G_n : \mathcal{A} \rightarrow \mathcal{A}$ that is similar to the zeroth-order gradient of f through n samples:

$$G_n(\mathbf{a}) := \frac{m}{n} \sum_{i=1}^n \left(T \left(\left((1 - \frac{1}{2p})\mathbf{a} + \frac{1}{2p}\mathbf{z}_i \right)^{\otimes p} \right) - T \left(\frac{1}{2p}\mathbf{z}_i \right) \right) \mathbf{z}_i + (\eta_i - \eta'_i)\mathbf{z}_i.$$

where $\mathbf{z}_i \sim \mathcal{N}(0, \frac{1}{m}\mathbf{I})$ and η_i, η'_i are independent zero-mean 1-sub-Gaussian noise. Therefore we have:

$$\begin{aligned} \mathbb{E}[G_n(\mathbf{a})] &= m \mathbb{E} \left[\sum_{l=0}^{p-1} \binom{p}{l} \mathbf{T} \left((1 - \frac{1}{2p})^{p-l} \mathbf{a}^{\otimes (p-l)} \otimes (\frac{1}{2p})^l \mathbf{z}^{\otimes l} \right) \mathbf{z} \right] \\ &\quad (\text{Due to symmetry of Gaussian only for odd } l =: 2s+1 \text{ expectation is nonzero}) \\ &= (1 - \frac{1}{2p})^{p-2s-1} (\frac{1}{2p})^{2s+1} \left[\sum_{s=0}^{\lfloor p/2-1 \rfloor} m^{-s} \binom{p}{2s+1} \mathbf{T}(\mathbf{a}^{\otimes (p-2s-1)} \otimes \mathbf{I}^{\otimes (s+1)}) \right] \end{aligned}$$

Note that for even p the last term (when $s = p/2 - 1$) is $\mathbf{T}(\mathbf{a} \otimes \mathbf{I}^{\otimes p/2}) = \sum_{j=1}^k \lambda_j (\mathbf{a}^\top \mathbf{v}_j) \mathbf{v}_j$. While all other terms will push the iterate towards the optimal action at a superlinear speed, the last term perform a matrix multiplication and the convergence speed will depend on the eigengap. Therefore for $p \geq 4$ we will remove the extra bias in the last term that is orthogonal to \mathbf{v}_1 and will treat it as noise. (Notice for quadratic function $s = 0 = p/2 - 1$ is the only term in $\mathbb{E}[G_n(\mathbf{a})]$. This is the distinction between $p = 2$ and larger p , and why its convergence depends on eigengap.)

We further define $G(\mathbf{a})$ as the population version of $G_n(\mathbf{a})$ by removing this undesirable bias term that will be treated as noise:

$$\begin{aligned} G(\mathbf{a}) &= \begin{cases} \mathbb{E}[G_n] - \frac{(\frac{1}{2p})^{p-1} (1 - \frac{1}{2p})^p}{m^{p/2-1}} \sum_{j=2}^k \lambda_j (\mathbf{v}_j^\top \mathbf{a}) \mathbf{v}_j, & \text{when } p \text{ is even} \\ \mathbb{E}[G_n], & \text{when } p \text{ is odd.} \end{cases} \\ &= \sum_{s=0}^{\lfloor (p-3)/2 \rfloor} \frac{(\frac{1}{2p})^{2s+1}}{m^s} \binom{p}{2s+1} \mathbf{T}(\mathbf{I}^{\otimes (s+1)} \otimes ((1 - \frac{1}{2p})\mathbf{a})^{\otimes (p-2s-1)}) \\ &= \frac{1}{2} (1 - \frac{1}{2p})^{p-1} \mathbf{T}(\mathbf{I}, \mathbf{a}^{\otimes (p-1)}) + O(1/m). \end{aligned}$$

We define $G(\mathbf{a})$ to push the action \mathbf{a} towards the \mathbf{v}_1 direction with at least linear convergence rate. More precisely, their angle $\tan \theta(G(\mathbf{a}), \mathbf{v}_1)$ will converge linearly to 0 for proper initialization with the dynamics $\mathbf{a} \rightarrow G(\mathbf{a})$. An easy way to see that is when $p = 2$ or 3, G is conducting (3-order tensor) power iteration. For higher-order problems, this operation G is equivalent to the summation of $p, p-2, p-4, \dots$ -th order tensor product and hence the linear convergence.

The estimation error $G_n(\mathbf{a}) - G(\mathbf{a})$ will be treated as noise (which is not mean zero when p is even but will be small enough: $O((2p)^{-p} m^{-(p-1)/2})$). Therefore the iterative algorithm with $\mathbf{a} \rightarrow G_n(\mathbf{a})$ will converge to a small neighborhood of \mathbf{v}_1 depending on the estimation error. This estimation error is controlled by the choice of sample size n in each iteration. We now provide the proof sketch:

Lemma D.3 (Initialization for $p \geq 3$; Corollary C.1 from [70]). For any $\eta \in (0, 1/2)$, with $L = \Theta(k \log(1/\eta))$ samples $\mathcal{A} = \{\mathbf{a}^{(1)}, \mathbf{a}^{(2)}, \dots, \mathbf{a}^{(L)}\}$ where each $\mathbf{a}^{(l)}$ is sampled uniformly on the sphere \mathbb{S}^{d-1} . At least one sample $\mathbf{a} \in \mathcal{A}$ satisfies

$$\max_{j \neq 1} |\mathbf{v}_j^\top \mathbf{a}| \leq 0.5 |\mathbf{v}_1^\top \mathbf{a}|, \text{ and } |\mathbf{v}_1^\top \mathbf{a}| \geq 1/\sqrt{d}. \quad (7)$$

with probability at least $1 - \eta$.

Lemma D.4 (Iterative progress). Let $\alpha = 1/2$ for $p \geq 3$ in Algorithm 6. Consider noisy operation $\mathbf{a}^+ \rightarrow G(\mathbf{a}) + \mathbf{g}$. If the error term \mathbf{g} satisfies:

$$\begin{aligned} \|\mathbf{g}\| &\leq \min\left\{\frac{0.025}{p} \lambda_1 (\mathbf{v}_1^\top \mathbf{a})^{p-2}, 0.1 \lambda_1 \tilde{\varepsilon}\right\} \\ &\quad + 0.03 \lambda_1 |\sin \theta(\mathbf{v}_1, \mathbf{a})| (\mathbf{v}_1^\top \mathbf{a})^{p-2}, \\ |\mathbf{v}_1^\top \mathbf{g}| &\leq 0.05 \lambda_1 (\mathbf{v}_1^\top \mathbf{a})^{p-1}. \end{aligned}$$

Suppose \mathbf{a} satisfies $0.5 |\mathbf{v}_1^\top \mathbf{a}| \geq \max_{j \geq 2} |\mathbf{v}_j^\top \mathbf{a}|$, we have:

$$\tan \theta(\mathbf{a}^+, \mathbf{v}_1) \leq 0.8 \tan \theta(\mathbf{a}, \mathbf{v}_1) + \tilde{\varepsilon}.$$

We can also bound \mathbf{g} by standard concentration plus an additional small bias term.

Lemma D.5 (Estimation error bound for G). For fixed value $\delta \in (0, 1)$ and large enough universal constant c_1, c_2, c_m, c_n , when $m = c_m d \log(n/\delta)$, $n \geq c_n d \log(d/\delta)$, we have

$$\begin{aligned} \|\mathbf{g}\| &\equiv \|G_n(\mathbf{a}) - G(\mathbf{a})\| \leq c_1 \sqrt{\frac{d^2 \log^3(n/\delta) \log(d/\delta)}{n}} + e \lambda_2 |\sin \theta(\mathbf{a}, \mathbf{v}_1)|, \\ |\mathbf{v}_1^\top \mathbf{g}| &\equiv |\mathbf{v}_1^\top G_n(\mathbf{a}) - \mathbf{v}_1^\top G(\mathbf{a})| \leq c_2 \sqrt{\frac{d \log^3(n/\delta) \log(d/\delta)}{n}}. \end{aligned}$$

with probability $1 - \delta$. $e = 0$ for odd p and $e = (2p)^{-(p-1)} m^{-(p/2-1)}$ for even p .

Together we are able to prove Theorem 3.14:

Proof of Theorem 3.14. Initially with high probability there exists an $\mathbf{a}_0 \in \mathcal{A}_0$ such that Eqn. (7) holds, i.e., $\mathbf{v}_1^\top \mathbf{a}_0 \geq 1/\sqrt{d}$ and $\mathbf{v}_j^\top \mathbf{a}_0 \leq 2|\mathbf{v}_1^\top \mathbf{a}_0|, \forall j \geq 2$.

Next, from Lemma D.5, the extra bias term is bounded by $e \lambda_2 |\sin \theta(\mathbf{a}, \mathbf{v}_1)| \leq 0.03 \lambda_1 (\mathbf{v}_1^\top \mathbf{a})^{p-2} |\sin \theta(\mathbf{a}, \mathbf{v}_1)|$ since $e = (2p)^{-p+1} m_s^{-p/2+1}$ and with our choice of variance $m_s \geq d \geq (\mathbf{v}_1^\top \mathbf{a})^{-2}$, plus $p \geq 3$. Next with our setting of $n_s = \tilde{\Theta}(d^p / (\lambda_1^2 \tilde{\varepsilon}_s^2))$, the error term $\|\mathbb{E}[G(\mathbf{a})] - G_n(\mathbf{a})\|$ is upper bounded by $\tilde{O}(\sqrt{\frac{d^2}{n}}) \leq 0.025 \lambda_1 d^{-(p-2)/2} \tilde{\varepsilon}_s / p + 0.1 \lambda_1 \tilde{\varepsilon}_s$. Meanwhile $|\mathbf{v}_1^\top \mathbf{g}| \leq \tilde{O}(\sqrt{\frac{d}{n_s}}) \leq 0.05 \lambda_1 (\mathbf{v}_1^\top \mathbf{a})^{p-1}$.

This meets the requirements for Theorem D.4 and therefore $\tan \theta(G_n(\mathbf{a}_0), \mathbf{v}_1) \leq 0.8 \tan \theta(\mathbf{a}_0, \mathbf{v}_1) + 0.1 \lambda_1 \tilde{\varepsilon}_s$. Therefore after l steps will have

$$\begin{aligned} \tan \theta(G_n^l(\mathbf{a}_0), \mathbf{v}_1) &\leq 0.8^l \tan \theta(\mathbf{a}_0, \mathbf{v}_1) + \sum_{i=1}^l 0.8^i \cdot 0.1 \tilde{\varepsilon}_s \\ &\leq 0.8^l \tan \theta(\mathbf{a}_0, \mathbf{v}_1) + 0.5 \tilde{\varepsilon}_s. \end{aligned}$$

Notice initially $\tan \theta(\mathbf{a}_0, \mathbf{v}_1) \leq 1/(\mathbf{v}_1^\top \mathbf{a}_0) \leq \sqrt{d}$. Therefore after at most $l = O(\log_2(\tan \theta(\mathbf{a}_0, \mathbf{v}_1))) \leq O(\log_2(d))$ steps, we will have $\tan \theta(G_n^l(\mathbf{a}_0), \mathbf{v}_1) \leq \tilde{\varepsilon}_0/2 = \tilde{\varepsilon}_1$. With the same argument, the progress also holds for $s > 0$ with even smaller l . \square

Proof of Lemma D.5. We first estimate $G_n(\mathbf{a}) - \mathbb{E}[G_n(\mathbf{a})]$, which is what we want for even p . For odd p we will need to analyze an extra bias term that is orthogonal to \mathbf{v}_1 , $\mathbf{e} := \frac{(\frac{1}{2p})^{p-1} (1 - \frac{1}{2p})^p}{m^{p/2-1}} \sum_{j=2}^k \lambda_j (\mathbf{v}_j^\top \mathbf{a}) \mathbf{v}_j$; and we have $G_n(\mathbf{a}) - \mathbb{E}[G_n(\mathbf{a})] = G_n(\mathbf{a}) - G(\mathbf{a}) + \mathbf{e}$.

We decompose $G_n(\mathbf{a})$ as $G_n(\mathbf{a}) = \sum_{s=1}^k G_n^{(s)} + N$, where $G_n^{(s)} := \frac{m}{n} \sum_{i=1}^n \binom{p}{s} T(((1 - 0.5/p)\mathbf{a})^{\otimes p-s} \otimes (\mathbf{z}_i/(2p))^{\otimes s}) \mathbf{z}_i$. The noise term $N := \frac{m}{n} \sum \epsilon_i \mathbf{z}_i$.

$$\begin{aligned} G_n^{(s)} &:= \frac{m}{n} \sum_{i=1}^n \binom{p}{s} T(((1 - 0.5/p)\mathbf{a})^{\otimes p-s} \otimes (\mathbf{z}_i/(2p))^{\otimes s}) \mathbf{z}_i \\ &= \frac{m}{n} \left(1 - \frac{1}{2p}\right)^{p-s} \left(\frac{1}{2p}\right)^s \binom{p}{s} \sum_{i=1}^n \sum_{j=1}^k \lambda_j (\mathbf{a}^\top \mathbf{v}_j)^{p-s} (\mathbf{z}_i^\top \mathbf{v}_j)^s \mathbf{z}_i. \\ \mathbb{E}[G_n^{(s)}] &= m \left(1 - \frac{1}{2p}\right)^{p-s} \left(\frac{1}{2p}\right)^s \binom{p}{s} \sum_{j=1}^k \lambda_j (\mathbf{a}^\top \mathbf{v}_j)^{p-s} \mathbb{E}[(\mathbf{z}^\top \mathbf{v}_j)^s \mathbf{z}] \\ &= \begin{cases} \left(1 - \frac{1}{2p}\right)^{p-s} \left(\frac{1}{2p}\right)^s m \binom{p}{s} \sum_{j=1}^k \lambda_j (\mathbf{a}^\top \mathbf{v}_j)^{p-s} \frac{1}{m^{(s+1)/2}} (s)!! \mathbf{v}_j, & \text{for odd } s, \\ 0, & \text{for even } s \end{cases} \\ &= \begin{cases} \left(1 - \frac{1}{2p}\right)^{p-s} \left(\frac{1}{2p}\right)^s \frac{s!!}{m^{(s-1)/2}} \binom{p}{s} \sum_{j=1}^k \lambda_j (\mathbf{a}^\top \mathbf{v}_j)^{p-s} \mathbf{v}_j, & \text{for odd } s, \\ 0, & \text{for even } s \end{cases} \end{aligned}$$

$$G_n^{(s)} - \mathbb{E}[G_n^{(s)}] = m \left(1 - \frac{1}{2p}\right)^{p-s} \left(\frac{1}{2p}\right)^s \binom{p}{s} \sum_{j=1}^k \lambda_j (\mathbf{a}^\top \mathbf{v}_j)^{p-s} \mathbf{g}_{n,s}(j),$$

where $\mathbf{g}_{n,s}(j) := \frac{1}{n} \sum_{i=1}^n (\mathbf{z}_i^\top \mathbf{v}_j)^s \mathbf{z}_i - \mathbb{E}[(\mathbf{z}^\top \mathbf{v}_j)^s \mathbf{z}]$.

Notice the scaling in each $G_n^{(s)}$ is $(1 - \frac{1}{2p})^{p-s} (\frac{1}{2p})^s \binom{p}{s} \leq (\frac{1}{2p})^s p^s / (s!) < 2^{-s}$ decays exponentially. In Claim D.12 we give bounds for $\mathbf{g}_{n,s}(j)$. We note the bound for each $\mathbf{g}_{n,s}$ also decays with s .

Therefore the bottleneck of the upper bound mostly depend on $\mathbf{g}_{n,0}$ and $\mathbf{v}_1^\top \mathbf{g}_{n,0}$, and we get:

$$\begin{aligned} \|G_n - \mathbb{E}[G_n]\| &\leq C_1 \lambda_1 \sqrt{\frac{(d + \log(1/\delta)) d \log(n/\delta)}{n}} + N, \\ |\mathbf{v}_1^\top G_n - \mathbb{E}[\mathbf{v}_1^\top G_n]| &\leq C_2 \lambda_1 \sqrt{\frac{d \log(n/\delta) (1 + \log(1/\delta))}{n}} + \mathbf{v}_1^\top N. \end{aligned}$$

Next from Claim D.11, the noise term

$$N \leq C_3 \sqrt{\frac{m \log(n/\delta) (d + \log(n/\delta)) \log(d/\delta)}{n}},$$

$$|\mathbf{v}_1^\top N| \leq C_4 \sqrt{m \frac{\log^2(n/\delta) \log(d/\delta)}{n}},$$

Finally \mathbf{e} is very small: $\|\mathbf{e}\| \leq \frac{1}{m^{p/2-1} (2p)^{(p-1)}} \lambda_2 \|\mathbf{V}^\top \mathbf{a}\| = \lambda_2 \frac{1}{m^{p/2-1} (2p)^{(p-1)}} \sin \theta(\mathbf{a}, \mathbf{v}_1)$. $|\mathbf{e}^\top \mathbf{v}_1| = 0$.

Together we can bound $G_n(\mathbf{a}) - G(\mathbf{a})$ and finish the proof. \square

Proof of Lemma D.4. From Lemma D.1 we have: $|\tan \theta(G(\mathbf{a}), \mathbf{v}_1)| \leq 1/2 |\tan \theta(\mathbf{a}, \mathbf{v}_1)|$. Let $\mathbf{V} = [\mathbf{v}_2, \dots, \mathbf{v}_k]$. For any $p \geq 2$, we have:

$$\begin{aligned} |\tan \theta(\mathbf{a}^+, \mathbf{v}_1)| &= \frac{\|\mathbf{V}^\top \mathbf{a}^+\|_2}{|\mathbf{v}_1^\top \mathbf{a}^+|} \\ &= \frac{\|\mathbf{V}^\top (G(\mathbf{a}) + \mathbf{g})\|}{|\mathbf{v}_1^\top (G(\mathbf{a}) + \mathbf{g})|} \\ &\leq \frac{\|\mathbf{V}^\top G(\mathbf{a})\| + \|\mathbf{V}^\top \mathbf{g}\|}{|\mathbf{v}_1^\top G(\mathbf{a})| - |\mathbf{v}_1^\top \mathbf{g}|} \\ &\leq \frac{1/2 |\tan \theta(\mathbf{a}, \mathbf{v}_1)| |\mathbf{v}_1^\top G(\mathbf{a})| + \|\mathbf{g}\|}{|\mathbf{v}_1^\top G(\mathbf{a})| - |\mathbf{v}_1^\top \mathbf{g}|} \\ &= \alpha |\tan \theta(\mathbf{a}, \mathbf{v}_1)| \frac{S_1}{S_1 - \|\mathbf{v}_1^\top \mathbf{g}\|} + \frac{\|\mathbf{g}\|}{S_1 - \|\mathbf{v}_1^\top \mathbf{g}\|}, \end{aligned}$$

where $S_1 := \mathbf{v}_1^\top G(\mathbf{a}) = \mathbf{v}_1^\top G(\mathbf{a}) = \lambda_1 \sum_{s=0}^{\lfloor (p-3)/2 \rfloor} \frac{(1-\frac{1}{2p})^{p-2s-1} (\frac{1}{2p})^{2s+1}}{m^s} \binom{p}{2s+1} (\mathbf{v}_1^\top \mathbf{a})^{p-2s-1} \geq \lambda_1 (1 - \frac{1}{2p})^{p-1} (\frac{1}{2p})^p (\mathbf{v}_1^\top \mathbf{a})^{p-1} \geq \frac{\lambda_1}{4} (\mathbf{v}_1^\top \mathbf{a})^{p-1}$. The inequality comes from keeping only the first term where $s = 0$. With the assumption that $|\mathbf{v}_1^\top \mathbf{g}| \leq 0.05 \lambda_1 (\mathbf{v}_1^\top \mathbf{a})^{p-1}$, we have $|\mathbf{v}_1^\top \mathbf{g}| \leq 0.2 S_1$. Therefore

$$\begin{aligned} |\tan \theta(\mathbf{a}^+, \mathbf{v}_1)| &\leq 1.25/2 |\tan \theta(\mathbf{a}, \mathbf{v}_1)| + \frac{\|\mathbf{g}\|}{S_1 - |\mathbf{v}_1^\top \mathbf{g}|} \\ &\leq 1.25/2 |\tan \theta(\mathbf{a}, \mathbf{v}_1)| + 5/4 \frac{\|\mathbf{g}\|}{S_1} \\ &\leq 1.25/2 |\tan \theta(\mathbf{a}, \mathbf{v}_1)| + 5 \frac{\|\mathbf{g}\|}{\lambda_1 (\mathbf{v}_1^\top \mathbf{a})^{p-1}}. \end{aligned}$$

Notice when $5 \frac{\|\mathbf{g}\|}{\lambda_1 (\mathbf{v}_1^\top \mathbf{a})^{p-1}} \leq \max\{0.125 |\tan \theta(\mathbf{a}, \mathbf{v}_1)|, \tilde{\varepsilon}\}$, which will ensure $|\tan \theta(G(\mathbf{a}), \mathbf{v}_1)| \leq (1.25/2 + 0.125) |\tan \theta(G(\mathbf{a}), \mathbf{v}_1)| + \tilde{\varepsilon}$. (We will prove this condition is satisfied when $\|\mathbf{g}\| \leq \min\{\frac{0.025}{p} \lambda_1 (\mathbf{v}_1^\top \mathbf{a})^{p-2}, 0.1 \lambda_1 \tilde{\varepsilon}\}$. We will handle the additional term in the upper bound of $\|\mathbf{g}\|$ later.) We divide this requirement into the following two cases. On one hand, when $|\mathbf{v}_1^\top \mathbf{a}| \leq 1 - 1/(p-1)$, $\|\mathbf{V}^\top \mathbf{a}\| \geq \sqrt{1 - (1 - 1/(p-1))^2} > 1/p$, therefore $|\tan \theta(\mathbf{a}, \mathbf{v}_1)| \geq 1/p |\mathbf{v}_1^\top \mathbf{a}|$. Therefore

$$\begin{aligned} 5 \frac{\|\mathbf{g}\|}{\lambda_1 (\mathbf{v}_1^\top \mathbf{a})^{p-1}} &\leq 0.125 |\tan \theta(\mathbf{a}, \mathbf{v}_1)| \\ \Leftrightarrow 5 \frac{\|\mathbf{g}\|}{\lambda_1 (\mathbf{v}_1^\top \mathbf{a})^{p-1}} &\leq 0.125 / (p |\mathbf{v}_1^\top \mathbf{a}|) \\ \Leftrightarrow \|\mathbf{g}\| &\leq 0.025 \lambda_1 (\mathbf{v}_1^\top \mathbf{a})^{p-2} / p. \end{aligned}$$

On the other hand, when $|\mathbf{v}_1^\top \mathbf{a}| \geq 1 - 1/(p-1)$, $|\mathbf{v}_1^\top \mathbf{a}|^{p-1} \geq 1/4$ when $p = 3$. Therefore $5 \frac{\|\mathbf{g}\|}{\lambda_1 (\mathbf{v}_1^\top \mathbf{a})^{p-1}} \leq 20 \|\mathbf{g}\| / \lambda_1$. Therefore we will need $\|\mathbf{g}\| \leq 0.05 \lambda_1 \tilde{\varepsilon}$, and then the requirement that $5 \frac{\|\mathbf{g}\|}{\lambda_1 (\mathbf{v}_1^\top \mathbf{a})^{p-1}} \leq \tilde{\varepsilon}$ is satisfied.

Altogether in both cases we have: $|\tan \theta(G(\mathbf{a}), \mathbf{v}_1)| \leq 0.75 |\tan \theta(G(\mathbf{a}), \mathbf{v}_1)| + \tilde{\varepsilon}$. Finally if we additionally increase $\|\mathbf{g}\|$ by $0.05 \lambda_1 (\mathbf{v}_1^\top \mathbf{a})^{p-1}$ we will have: $|\tan \theta(G(\mathbf{a}), \mathbf{v}_1)| \leq 0.8 |\tan \theta(G(\mathbf{a}), \mathbf{v}_1)| + \tilde{\varepsilon}$. \square

Proof of Corollary 3.16. As shown in Theorem 3.14 at least one action \mathbf{a} in \mathcal{A}_S , $|\mathcal{A}_S| \leq \tilde{O}(k)$ satisfies $\tan \theta(\mathbf{a}, \mathbf{a}^*) \leq \varepsilon$ with a total of $\tilde{O}(\frac{d^p k}{\lambda_1^2 \varepsilon^2})$ steps. Therefore with Claim D.6 we have to get $\tilde{\varepsilon}$ -optimal reward we need $\tilde{O}(\frac{d^p k}{\lambda_1 \tilde{\varepsilon}})$ steps. Notice the eluder dimension for symmetric polynomials is d^p and the size of \mathcal{A}_S is at most $\tilde{O}(k)$. Then by applying Corollary D.9 we get that the total regret is at most $\tilde{O}(\sqrt{d^p k T} + \sqrt{|\mathcal{A}_S| T}) = \tilde{O}(\sqrt{d^p k T})$. \square

D.2 PAC to Regret Bound Relation.

Claim D.6 (Connecting angle to regret). *When $0 < \tan \theta(\mathbf{a}, \mathbf{v}_1) \leq \zeta$, we have regret $r^* - r(\mathbf{a}) \leq r^* \min\{2, p\zeta^2\}$.*

Proof.

$$\begin{aligned}
|\cos \theta(\mathbf{a}, \mathbf{v}_1)| &= |\mathbf{a}^\top \mathbf{v}_1| =: b, \\
|\tan \theta(\mathbf{a}, \mathbf{v}_1)| &= \frac{\sqrt{1-b^2}}{b} \leq \zeta \Leftrightarrow b \geq \frac{1}{\sqrt{\zeta^2+1}}. \\
\Rightarrow r^* - r(\mathbf{a}) &\leq \lambda_1 - \lambda_1 b^p \\
&\leq \lambda_1 (1 - (\zeta^2 + 1)^{-p/2}) \\
&= \lambda_1 \frac{(\zeta^2 + 1)^{p/2} - 1}{(\zeta^2 + 1)^{p/2}} \\
&\leq \lambda_1 ((\zeta^2 + 1)^{p/2} - 1) \quad (\text{since denominator } (\zeta^2 + 1)^{p/2} \geq 1) \\
&\leq \lambda_1 p \zeta^2, \text{ when } \zeta^2 \leq 1/p.
\end{aligned}$$

Additionally by definition $r^* - r(\mathbf{a}) \leq \lambda_1 - (-\lambda_1) = 2\lambda_1$ and thus $r^* - r(\mathbf{a}) \leq \lambda_1 \min\{2, p\zeta^2\}$. We now derive the last inequality. When $\zeta \geq 1/p$ it is trivially true. When $\zeta \leq 1/p$, we have $(1 + \zeta^2)^{p/2} \leq 1 + p\zeta^2$ for any $p \geq 2$. Since the LHS is a convex function for ζ when $p \geq 2$ and when $\zeta = 0$ LHS=RHS and when $\zeta^2 = 1/p$ LHS is always smaller than RHS (=2).

Notice the argument is straightforward to extend to the setting where the angle is between \mathbf{a} and subspace V_1 that satisfies $\forall \mathbf{v} \in V_1, T(\mathbf{v}) \geq \lambda_1 - \epsilon$, then one also get $r^* - r(\mathbf{a}) \leq \lambda_1 - (\lambda_1 - \epsilon)b^p \leq \min\{\lambda_1, \lambda_1 p \zeta^2 + \epsilon b^p\} \leq \min\{\lambda_1, \lambda_1 p \zeta^2 + \epsilon\}$. \square

Claim D.7 (Connecting PAC to Cumulative Regret). *Suppose we have an algorithm $\text{alg}(\zeta)$ that finds ζ -optimal action $\hat{\mathbf{a}}$ that satisfies $0 < \tan \theta(\mathbf{a}, \mathbf{v}_1) \leq \zeta$ by taking $A\zeta^{-a}$ actions. Here A can depend on any parameters such as $d, \lambda_1, \text{probability error } \delta, \text{ etc.},$ that are not ζ . Then for large enough T , by calling alg with $\zeta = A^{\frac{1}{a+2}} T^{-\frac{1}{a+2}} p^{-\frac{1}{a+2}}$ and playing its output action $\hat{\mathbf{a}}$ for the remaining actions, one can get a cumulative regret of:*

$$\mathfrak{R}(T) \lesssim T^{\frac{a}{a+2}} p^{\frac{a}{a+2}} A^{\frac{2}{a+2}} r^*.$$

Similarly, if an oracle finds ϵ -optimal action $\hat{\mathbf{a}}$ that satisfies $r^* - r(\mathbf{a}) \leq \epsilon$ with $B\epsilon^{-b}$ samples, then by setting $\epsilon = (Br^*/T)^{\frac{1}{1+b}}$, and playing the output arm for the remaining actions, one can get cumulative regret of:

$$\mathfrak{R}(T) \lesssim B^{\frac{1}{1+b}} T^{\frac{b}{1+b}} r^{\frac{1}{1+b}}.$$

Proof. For the chosen ζ , write $T_1 = A\zeta^{-a}$ be the number of actions that finds ζ -optimal action. Therefore $T_1 = A^{\frac{2}{a+2}} T^{\frac{a}{a+2}} p^{\frac{a}{a+2}}$. First, when $T \geq Ap^{a/2}$, $\zeta^2 \leq 1/p$, namely $r^* - r(\mathbf{a}) \leq r^* p \zeta^2$. We have:

$$\begin{aligned}
\mathfrak{R}(T) &\leq \sum_{t=1}^{T_1} 2r^* + \sum_{t=T_1+1}^T r^* p \zeta^2 \\
&\leq 2r^* T_1 + T r^* p \zeta^2 \\
&\leq 3T^{\frac{a}{a+2}} p^{\frac{a}{a+2}} A^{\frac{2}{a+2}} r^*.
\end{aligned}$$

When $T < Ap^{a/2}$, it trivially holds that $\mathfrak{R}(T) \leq 2r^* T < 2T^{\frac{a}{a+2}} p^{\frac{a}{a+2}} A^{\frac{2}{a+2}} r^*$. \square

Theorem D.8 (Theorem 5.1 from [5]). *With UCB algorithm on action set with size K , we have with probability $1 - \delta$,*

$$\mathfrak{R}(T) = \tilde{O}(\min\{\sqrt{KT}\} + K).$$

Algorithm 7 UCB (Algorithm 1 in Section 5 of [5])

- 1: **Input:** Stochastic reward function f , failure probability δ , action set \mathcal{A} with finite size K .
 - 2: **for** t from 1 to $T - 1 - K$ **do**
 - 3: Execute arm $I_t = \arg \max_{i \in [K]} \left(\hat{\mu}^t(i) + \sqrt{\frac{\log(TK/\delta)}{N^t(i)}} \right)$. Here $N^t(\mathbf{a}) = 1 + \sum_{i=1}^t \mathbf{1}\{I_i = \mathbf{a}\}$; and $\hat{\mu}^t(\mathbf{a}) = \frac{1}{N^t(\mathbf{a})} \left(r_{\mathbf{a}} + \sum_{i=1}^t \mathbf{1}\{I_i = \mathbf{a}\} r_i \right)$.
 - 4: Observe r_{I_t}
-

Corollary D.9. *With the same setting of Claim D.7, except that now the algorithm $\text{alg}(\varepsilon)$ finds a set \mathcal{A} of size S where at least one action $\mathbf{a} \in \mathcal{A}$ satisfies $r^* - f(\mathbf{a}) \leq \varepsilon$. Then all argument in Claim D.7 still hold by adding $\tilde{O}(\sqrt{ST})$ on the RHS of each regret bound.*

Proof. Suppose we run alg for T_1 steps and achieve ε -optimal reward.

Let $r_\varepsilon := \max_{\mathbf{a} \in \mathcal{A}} f(\mathbf{a})$. Therefore with UCB on mutiarm bandit we have: $\sum_{t=T_1+1}^T r_\varepsilon - f(\mathbf{a}_t) \leq \tilde{O}(\sqrt{ST})$ by Theorem D.8.

From the statement $r_\varepsilon \geq r^* - \zeta$. Therefore $\sum_{t=T_1+1}^T r^* - f(\mathbf{a}_t) \leq \tilde{O}(\sqrt{ST}) + \varepsilon(T - T_1)$. Therefore

$$\mathfrak{R}(T) \leq \sum_{t=1}^{T_1} 2r^* + \varepsilon(T - T_1) + \tilde{O}(\sqrt{ST}).$$

With the same choices of T_1 in Claim D.7, the same conclusion still holds with an additional term of $\tilde{O}(\sqrt{ST})$. \square

For symmetric tensor problems the set size is $\tilde{O}(k)$ and therefore we will have an additional \sqrt{kT} term which will be subsumed in our regret bound.

D.3 Variance and Noise Concentration

Lemma D.10 (Vector Bernstein; adapted from Theorem 7.3.1 in [65]). *Consider a finite sequence $\{\mathbf{x}_k\}_{k=1}^n$ be i.i.d randomly generated samples, $x_k \in \mathbb{R}^d$, and assume that $\mathbb{E}[\mathbf{x}_k] = 0$, $\|\mathbf{x}_k\| \leq L$, and covariance matrix of x_k is Σ . Then it satisfies that when $n \geq \log d/\delta$, we have:*

$$\left\| \frac{\sum_{i=1}^n \mathbf{x}_i}{n} \right\| \leq C \sqrt{\frac{(\|\Sigma\| + L^2) \log d/\delta}{n}},$$

with probability $1 - \delta$.

Claim D.11 (Noise concentration). *Let independent samples $\mathbf{z}_i \sim \mathcal{N}(0, 1/mI_d)$ and $\epsilon_i \sim \mathcal{N}(0, 1)$. With probability $1 - \delta$, $\delta \in (0, 1)$:*

$$\left\| \frac{m}{n} \sum_{i=1}^n \epsilon_i \mathbf{z}_i \right\| \leq C \sqrt{\frac{m \log(n/\delta) (d + \log(n/\delta)) \log(d/\delta)}{n}}$$

$$\left| \frac{m}{n} \sum_{i=1}^n \epsilon_i \mathbf{z}_i^\top \mathbf{v}_1 \right| \leq C' \sqrt{\frac{m \log^2(n/\delta) \log(d/\delta)}{n}}.$$

Proof. We use the Vector Bernstein Lemma D.10. The covariance matrix for $\mathbf{x}_i = \epsilon_i \mathbf{z}_i$ satisfies $\mathbb{E}[\mathbf{x}_i \mathbf{x}_i^\top] = 1/mI_d$. \mathbf{x}_i is mean zero. $\|\epsilon_i \mathbf{z}_i\|^2 = \epsilon_i^2 \|\mathbf{z}_i\|^2$. Notice $\epsilon_i^2 \sim \chi(1) \lesssim 1 + \log(1/\delta)$ and $m \mathbf{z}_i^\top \mathbf{z}_i \sim \chi(d) \lesssim d + \log(1/\delta)$. Therefore by directly applying Vector Bernstein $\|\epsilon_i \mathbf{z}_i\| \leq c \sqrt{\frac{(1 + \log(1/\delta))(d + \log(1/\delta))}{m}}$ with probability $1 - \delta$. By union bound we have: for all i , $\|\epsilon_i \mathbf{z}_i\| \leq c \sqrt{\frac{\log(n/\delta)(d + \log(n/\delta))}{m}}$ with probability $1 - \delta$. Therefore

$$\left\| \frac{1}{n} \sum_{i=1}^n \epsilon_i \mathbf{z}_i \right\| \leq C \sqrt{\frac{\log(n/\delta)(d + \log(n/\delta)) \log(d/\delta)}{mn}},$$

with probability $1 - \delta$. Similarly

$$\begin{aligned} \left| \frac{1}{n} \sum_{i=1}^n \epsilon_i \mathbf{z}_i^\top \mathbf{v}_1 \right| &\leq C \sqrt{\frac{\log(n/\delta)(1 + \log(n/\delta)) \log(d/\delta)}{mn}} \\ &= C' \sqrt{\frac{\log^2(n/\delta) \log(d/\delta)}{mn}}, \end{aligned}$$

□

Claim D.12. Let $\{\mathbf{z}_i\}_{i=1}^n$ be i.i.d samples from $\mathcal{N}(0, 1/mI_d)$. Let $g_{n,s}(j) := \frac{1}{n} \sum_{i=1}^n (\mathbf{z}_i^\top \mathbf{v}_j)^s \mathbf{z}_i - \mathbb{E}[(\mathbf{z}^\top \mathbf{v}_j)^s \mathbf{z}]$. We have:

$$\begin{aligned} \|g_{n,0}(j)\| &\lesssim \sqrt{\frac{d + \log(1/\delta)}{nm}}, \\ |\mathbf{v}_1^\top g_{n,0}(j)| &\lesssim \sqrt{\frac{1 + \log(1/\delta)}{nm}}, \\ |\mathbf{v}_1^\top g_{n,1}(j)| \leq \|g_{n,1}(j)\| &\lesssim \sqrt{\frac{d + \log(1/\delta)}{m^2 n}}, \text{ when } n \geq d \log(1/\delta), \\ |\mathbf{v}_1^\top g_{n,s}(j)| \leq \|g_{n,s}(j)\| &\lesssim \sqrt{\frac{\log(d/\delta)}{d^s n}}, \text{ when } n \geq \log(d/\delta), m \geq c_0 d \log(n/\delta), s \geq 2. \end{aligned}$$

For any $j \in [k]$.

We mostly care about the correct concentration for smaller s . For larger s a very loose bound will already suffice our requirement.

Proof of Claim D.12. For $s = 0$, $nm \|\frac{1}{n} \sum_{i=1}^n \mathbf{z}_i\|^2 \sim \chi(d)$, therefore $\|\frac{1}{n} \sum_{i=1}^n \mathbf{z}_i\| \lesssim \sqrt{\frac{d + \log(1/\delta)}{nm}}$. $nm(\frac{1}{n} \sum_{i=1}^n \mathbf{z}_i^\top \mathbf{v}_1)^2 \sim \chi(1)$. Therefore $|\frac{1}{n} \sum_{i=1}^n \mathbf{z}_i^\top \mathbf{v}_1| \lesssim \sqrt{\frac{1 + \log(1/\delta)}{nm}}$.

For $s = 1$, due to standard concentration for covariance matrices (see e.g. [65, 22]), we have:

$$m \left\| \left(\frac{1}{n} \sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i^\top - \mathbb{E}[\mathbf{z} \mathbf{z}^\top] \right) \right\| \leq \max \left\{ \sqrt{\frac{d + \log(2/\delta)}{n}}, \frac{d + \log(2/\delta)}{n} \right\}.$$

Therefore when $n \geq d \log(1/\delta)$, both results

$$\begin{aligned} \|g_{n,1}(j)\| &\lesssim \sqrt{\frac{d + \log(1/\delta)}{m^2 n}} \|\mathbf{v}_j\|, \\ &= \sqrt{\frac{d + \log(1/\delta)}{m^2 n}}, \text{ and} \\ \|\mathbf{v}_1^\top g_{n,1}(j)\| &\lesssim \sqrt{\frac{d + \log(1/\delta)}{m^2 n}} \|\mathbf{v}_1\| \|\mathbf{v}_j\| \\ &= \sqrt{\frac{d + \log(1/\delta)}{m^2 n}} \end{aligned}$$

hold.

For larger $s \geq 2$, with probability $1 - \delta$, $|\mathbf{z}_i^\top \mathbf{v}_j| \leq C \sqrt{\log(n/\delta)/m} = C c_0 / \sqrt{d} \leq 1/\sqrt{d}$. When $m \geq c_0 d \log(n/\delta)$, for small enough c_0 we have $|\mathbf{z}_i^\top \mathbf{v}_j| \leq 1/\sqrt{d}$ and $\|\mathbf{z}_i\| \leq 1$ for all $i \in [n]$. Therefore $\|(\mathbf{z}_i^\top \mathbf{v}_j)^s \mathbf{z}_i\| \leq d^{-s/2}$. We can use vector Bernstein, i.e., Lemma D.10 to get:

$$\|g_{n,s}(j)\| \leq C_1 \sqrt{\frac{\log(d/\delta)}{d^s n}}.$$

Therefore we have:

$$|g_{n,s}(j)^\top \mathbf{v}_1| \leq C_1 \sqrt{\frac{\log(d/\delta)}{d^s n}}.$$

□

D.3.1 The asymmetric setting

Now we consider the asymmetric tensor problem with reward $f : \mathcal{A} \rightarrow \mathbb{R}$. The input space \mathcal{A} consists of p vectors in a unit ball: $\vec{\mathbf{a}} = (\mathbf{a}(1), \mathbf{a}(2), \dots, \mathbf{a}(p)) \in \mathcal{A}, \|\mathbf{a}(s)\| \leq 1, \forall s \in [p]$. $f(\vec{\mathbf{a}}) = \mathbf{T}(\otimes_{s=1}^p \mathbf{a}(s)) + \eta$. Tensor $\mathbf{T} = \sum_{j=1}^k \lambda_j \mathbf{v}_j(1) \otimes \mathbf{v}_j(2) \cdots \otimes \mathbf{v}_j(p)$. For each $s \in [p]$, $\{\mathbf{v}_1(s), \mathbf{v}_2(s), \dots, \mathbf{v}_k(s)\}$ are orthonormal vectors. We order the eigenvalues such that $\lambda_1 \geq |\lambda_2| \cdots \geq |\lambda_k|$. Therefore the optimal reward is λ_1 and can be achieved by $\mathbf{a}^*(s) = \mathbf{v}_1(s), s \in [p]$. In this section we only consider $p \geq 3$ and leave the quadratic and low-rank matrix setting to the next section.

Theorem D.13. *For $p \geq 3$, by conducting alternating power iteration, one can get a ε -optimal reward with a total $\tilde{O}((2k)^p \log^p(p/\delta) d^p \lambda_1^{-1} \varepsilon^{-1})$ actions; therefore the regret bound is at most $\tilde{O}(\sqrt{k^p d^p T})$.*

This setting is actually much easier than the symmetric setting. Notice by replacing one slice of $\vec{\mathbf{a}}$ by random Gaussian $\mathbf{z}_i \sim \mathcal{N}(0, 2/d \log(d/\delta))$, one directly gets $\mathbf{T}(\mathbf{a}(1), \dots, \mathbf{a}(s-1), \mathbf{I}, \mathbf{a}(s+1), \dots, \mathbf{a}(p))$ on each slice with $1/n \sum_i f(\mathbf{a}(1), \dots, \mathbf{a}(s-1), \mathbf{z}_i, \mathbf{a}(s+1), \dots, \mathbf{a}(p)) \mathbf{z}_i$ which is tensor product. We defer the proof to Appendix D.4.

D.4 Omitted Details for Asymmetric Tensors

Algorithm 8 Phased elimination with alternating tensor product.

- 1: **Input:** Stochastic reward $r : (B_1^d)^{\otimes p} \rightarrow \mathbb{R}$ of polynomial degree p , failure probability δ , error ε .
 - 2: **Initialization:** $L_0 = C_L k \log(1/\delta)$; Total number of stages $S = C_S \lceil \log(1/\varepsilon) \rceil + 1$, $\mathcal{A}_0 = \{\mathbf{a}_0^{(1)}, \mathbf{a}_0^{(2)}, \dots, \mathbf{a}_0^{(L_0)}\} \subset (B_1^d)^{\otimes p}$ where each $\mathbf{a}_0^{(l)}(j), j \in [p]$ is uniformly sampled on the unit sphere \mathbb{S}^{d-1} . $\tilde{\varepsilon}_0 = 1$.
 - 3: **for** s from 1 to S **do**
 - 4: $\tilde{\varepsilon}_s \leftarrow \tilde{\varepsilon}_{s-1}/2$, $n_s \leftarrow C_n d^p \log(d/\delta) / \tilde{\varepsilon}_s^2$, $n_s \leftarrow n_s \cdot \log^3(n_s/\delta)$, $m_s \leftarrow C_m d \log(n/\delta)$, $\mathcal{A}_s = \emptyset$.
 - 5: **for** l from 1 to L_{s-1} **do**
 - 6: **Tensor product update:**
 - 7: Locate current arm $\tilde{\mathbf{a}} = \mathbf{a}_{s-1}^{(l)}$.
 - 8: **for** $\lceil (\lambda_1/\Delta) \log(2d) \rceil$ times **do**
 - 9: **for** j from 1 to p **do**
 - 10: Sample $\mathbf{z}_i \sim \mathcal{N}(0, 1/m_s I_d), i = 1, 2, \dots, n_s$.
 - 11: Calculate tentative arm $\mathbf{a}_i \leftarrow \tilde{\mathbf{a}}, \mathbf{a}_i(j) = (1 - \tilde{\varepsilon}_s) \tilde{\mathbf{a}}(j) + \tilde{\varepsilon}_s \mathbf{z}_i$
 - 12: Conduct estimation $\mathbf{y} \leftarrow 1/n_s \sum_{i=1}^{n_s} r_{\varepsilon_i}(\mathbf{a}_i) \mathbf{z}_i$.
 - 13: Update the current arm $\tilde{\mathbf{a}}(j) \leftarrow \mathbf{y} / \|\mathbf{y}\|$.
 - 14: Estimate the expected reward for $\tilde{\mathbf{a}}$ through n_s samples: $r_n = 1/n_s \sum_{i=1}^{n_s} r_{\varepsilon_i}(\tilde{\mathbf{a}})$.
 - 15: **Candidate Elimination:**
 - 16: **if** $r_n \geq \lambda_1(1 - p\tilde{\varepsilon}_s)$ **then**
 - 17: Keep the arm $\mathcal{A}_s \leftarrow \mathcal{A}_s \cup \{\tilde{\mathbf{a}}\}$
 - 18: Label the arms: $L_s = |\mathcal{A}_s|, \mathcal{A}_t = \{\mathbf{a}_s^{(1)}, \dots, \mathbf{a}_s^{(L_s)}\}$.
 - 19: Run Algorithm 7 with \mathcal{A}_S .
-

Lemma D.14 (Asymmetric Tensor Initialization). *With probability $1 - \delta$, with $L = \tilde{\Theta}((2k)^p \log^p(p/\delta))$ random initializations $\mathcal{A}_0 = \{\mathbf{a}_0^{(0)}, \mathbf{a}_0^{(1)}, \dots, \mathbf{a}_0^{(L)}\}$, there exists an initialization $\mathbf{a}_0 \in \mathcal{A}_0$ that satisfies:*

$$\begin{aligned} \alpha \mathbf{a}_0(s)^\top \mathbf{v}_1^{(s)} &\geq |\mathbf{a}_0(s)^\top \mathbf{v}_j^{(s)}|, \forall j \geq 2 \& j \in [k], \forall s \in [p], \\ \mathbf{a}_0(s)^\top \mathbf{v}_1^{(s)} &\geq 1/\sqrt{d}. \end{aligned} \tag{8}$$

with some constant $\alpha < 1$.

Proof. This lemma simply comes from applying Lemma D.3 for p times and we need $\geq 2k \log_2(p\delta)$ to ensure the condition for each $\mathbf{a}_0(s), s \in [p]$ holds. Therefore together we will need $(2k \log_2(p/\delta))^p$ samples. \square

Lemma D.15 (Asymmetric tensor progress). *For each a that satisfies Eqn. (8) with constant $\alpha < 1$, we have:*

$$\tan \theta(\mathbf{T}(\mathbf{a}(1), \dots, \mathbf{a}(s-1), \mathbf{I}, \mathbf{a}(s+1), \dots, \mathbf{a}(p)), \mathbf{v}_1^{(s)}) \leq \alpha \tan \theta(\mathbf{a}_j, \mathbf{v}_1^{(j)}),$$

for any j that is in $[p]$ but is not s . When $n \geq \Theta(d^p \log(d/\delta) \log^3(n/\delta)/\tilde{\varepsilon}^2)$ and $m = \Theta(d \log(n/\delta))$, we have:

$$\begin{aligned} & \tan \theta(\mathbf{T}(\mathbf{a}(1), \dots, \mathbf{a}(s-1), \mathbf{I}, \mathbf{a}(s+1), \dots, \mathbf{a}(p)), \mathbf{v}_1^{(s)}) \\ & \leq (1 + \alpha)/2 \tan \theta(\mathbf{a}_j, \mathbf{v}_1^{(j)}) + \tilde{\varepsilon}, \forall j \in [p] \& j \neq s. \end{aligned}$$

The remaining proof is a simpler version for the symmetric tensor setting on conducting noisy power method with the good initialization and iterative progress.

Finally due to the good initialization that satisfies (8) and together with Lemma D.15 we can finish the proof for Theorem D.13.

E Proof of Theorem 3.21

E.1 Additional Notations

Here, we briefly introduce complex and real algebraic geometry. This section is based on [56, 64, 12, 69].

An (affine) **algebraic variety** is the common zero loci of a set of polynomials, defined as $V = Z(S) = \{\mathbf{x} \in \mathbb{C}^n : f(\mathbf{x}) = 0, \forall f \in S\} \subseteq \mathbb{A}^n = \mathbb{C}^n$ for some $S \subseteq \mathbb{C}[x_1, \dots, x_n]$. A **projective variety** U is a subset of $\mathbb{P}^n = (\mathbb{C}^{n+1} \setminus \{0\})/\sim$, where $(x_0, \dots, x_n) \sim k(x_0, \dots, x_n)$ for $k \neq 0$ and S is a set of homogeneous polynomials of $(n+1)$ variables.

For an affine variety V , its **projectivization** is the variety $\mathbb{P}(V) = \{[\mathbf{x}] : \mathbf{x} \in V\} \subseteq \mathbb{P}^{n-1}$, where $[\mathbf{x}]$ is the line corresponding to \mathbf{x} .

The **Zariski topology** is the topology generated by taking all varieties to be the closed sets.

A set is **irreducible** if it is not the union of two proper closed subsets.

A **variety is irreducible** if and only if it is irreducible under the Zariski topology.

The **algebraic dimension** $d = \dim V$ of a variety V is defined as the length of the longest chain $V_0 \subset V_1 \subset \dots \subset V_d = V$, such that each V_i is irreducible.

A variety V is said to be **admissible** to a set of linear functions $\{\ell_\alpha : \mathbb{C}^d \rightarrow \mathbb{C}\}_{\alpha \in I}$, if for every ℓ_α , we have $\dim(V \cap \{\mathbf{x} \in \mathbb{C}^d : \ell_\alpha(\mathbf{x}) = 0\}) < \dim V$.

A map $f = (f_1, \dots, f_m) : \mathbb{A}^n \rightarrow \mathbb{A}^m$ is **regular** if each f_i is a polynomial.

A **algebraic set** is the common real zero loci of a set of polynomials.

For a complex variety $V \subseteq \mathbb{A}^n$, its real points form an algebraic set $V_{\mathbb{R}}$.

For an algebraic set $V_{\mathbb{R}}$, its real dimension $d = \dim_{\mathbb{R}} V_{\mathbb{R}}$ is the maximum number d such that $V_{\mathbb{R}}$ is locally semi-algebraically homeomorphic to the unit cube $(0, 1)^d$, details can be found in [12].

E.2 Proof of Sample Complexity

Lemma E.1 ([69], Theorem 3.2). *For $i = 1, \dots, T$, let $L_i : \mathbb{C}^n \times \mathbb{C}^m \rightarrow \mathbb{C}$ be bilinear functions and V_i be varieties given by homogeneous polynomials in \mathbb{C}^n . Let $V = V_1 \times \dots \times V_T \subseteq (\mathbb{C}^n)^N$. Let $W \subseteq \mathbb{C}^m$ be a variety given by homogeneous polynomials. In addition, we assume V_i is admissible with respect to the linear functions $\{f^{\mathbf{w}}(\cdot) = L_i(\cdot, \mathbf{w}) : \mathbf{w} \in W \setminus \{0\}\}$. When $T \geq \dim W$, let $\delta = T - \dim W + 1 \geq 1$. Then there exists a subvariety $Z \subseteq V$ with $\dim Z \leq \dim V - \delta$ such that for any $(\mathbf{x}_1, \dots, \mathbf{x}_T) \in V \setminus Z$ and $\mathbf{w} \in W$, if $L_1(\mathbf{x}_1, \mathbf{w}) = \dots = L_T(\mathbf{x}_T, \mathbf{w}) = 0$, then $\mathbf{w} = 0$.*

Lemma E.2 ([69], Lemma 3.1). *Let V be an algebraic variety in \mathbb{C}^d . Then $\dim_{\mathbb{R}} V_{\mathbb{R}} \leq \dim V$.*

Lemma E.3. *Let W be a vector space. For vectors $\mathbf{x}_1, \dots, \mathbf{x}_T$, if the map $f : \mathbf{w} \mapsto (\langle \mathbf{x}_1, \mathbf{w} \rangle, \dots, \langle \mathbf{x}_T, \mathbf{w} \rangle)$ is not injective over $W - W := \{\mathbf{w}_1 - \mathbf{w}_2 : \mathbf{w}_1, \mathbf{w}_2 \in W\}$, then there exists $\mathbf{v} \in W$ such that $f(\mathbf{v}) = 0$.*

Proof. Suppose $f(\mathbf{w}_1) = f(\mathbf{w}_2)$. Let $\mathbf{v} = \mathbf{w}_1 - \mathbf{w}_2$. Then $\mathbf{v} \in W - W$ and $f(\mathbf{v}) = f(\mathbf{w}_1) - f(\mathbf{w}_2) = 0$. \square

Definition E.4 (Tensorization). Let f be a polynomial of x_1, \dots, x_d with degree $\deg f \leq p$. Then every p -tensor \mathbf{W}_f satisfying $\langle \mathbf{W}_f, \mathbf{X}_x \rangle = f(\mathbf{x})$ is said to be a *tensorization* of the polynomial f , where \mathbf{X}_x is the tensorization of \mathbf{x} itself:

$$\mathbf{X}_x = \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix}^{\otimes p}. \quad (9)$$

Let \mathcal{F} be a class of polynomials. A variety of tensorization of \mathcal{F} is defined to be an irreducible closed variety defined by homogeneous polynomials W , such that for every $f \in \mathcal{F}$, there is a tensorization \mathbf{W}_f of f , such that $W \ni \mathbf{W}_f$ contains its tensorization. Note that neither tensorization of f nor variety of tensorization of \mathcal{F} is unique.

We define the variety of tensorization of \mathbf{x} as follows. (Note that this is uniquely defined.) Consider the regular map

$$\varphi_1 : \mathbb{C}^d \rightarrow \mathbb{C}^{(d+1)^p}, \quad \mathbf{x} \mapsto \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix}^{\otimes p}, \quad (10)$$

the tensorization of \mathbf{x} is defined as $V_i = \mathbb{P}(\overline{\text{Im } \varphi_1})$.

Note that V_i is irreducible because φ_1 is regular and \mathbb{C}^d is irreducible. By [56, Theorem 9.9], its dimension is given by

$$\dim V_i \leq \dim \overline{\text{Im } \varphi_1} + 1 \leq \dim \mathbb{C}^d + 1 = d + 1. \quad (11)$$

Lemma E.5. For any non-zero polynomial $f \neq 0$ with $\deg f \leq p$. Let W_f be a tensorization of f . Then V_i is admissible with respect to $\{L_i(\cdot) = \langle \cdot, W_f \rangle\}$.

Proof. Since V_i is irreducible and L_i is a linear function, it suffices to verify that $\langle \mathbf{X}_x, W_f \rangle \neq 0$ [69]. But according to Definition E.4, $\langle \mathbf{X}_x, W_f \rangle \neq 0$ is equivalent to

$$f(\mathbf{x}) = \left\langle \mathbf{W}_f, \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix}^{\otimes p} \right\rangle \neq 0. \quad (12)$$

Since $f \neq 0$, we must have $f(\mathbf{x}) \neq 0$ for some \mathbf{x} , which gives a non-zero $\mathbf{X}_x \neq 0$ for the above equation: $\langle \mathbf{X}_x, W_f \rangle \neq 0$, and we conclude that V_i is admissible. \square

Lemma E.6. Let $V \subset \mathbb{C}^n$ be a (Zariski) closed proper subset, $V \neq \mathbb{C}^n$. Then V is a null set, i.e. it has (Lebesgue) measure zero.

Proof. Suppose $V = Z(S)$ is the vanishing set for some $S \subseteq \mathbb{C}[x_1, \dots, x_n]$. Since $V \neq \mathbb{C}^n$, let $f \in S$, we have $V \subseteq Z(f)$, so it suffices to show $\text{Leb}(Z(f)) = 0$, which is because $Z(f) = f^{-1}(0)$, $\text{Leb}(\{0\}) = 0$, f is a continuous function (under Euclidean topology), and $\text{Leb}(\{\mathbf{x} : \nabla f(\mathbf{x}) = 0\}) = 0$. \square

Theorem E.7. Assume that the reward function class is a class of polynomials \mathcal{F} . Let W be (one of) its variety of tensorization. If we sample $T \geq \dim W$ times, and the sample points satisfying $(\mathbf{x}_1, \dots, \mathbf{x}_T) \in (\mathbb{C}^d)^T \setminus Z$ for some null set Z . Then we can uniquely determine the reward function f from the observed rewards $(f(\mathbf{x}_1), \dots, f(\mathbf{x}_T))$.

Proof. Let $n = m = (d+1)^p$, $L_i(\mathbf{x}, \mathbf{w}) = \langle \mathbf{x}, \mathbf{w} \rangle$, $V = V_1 \times \dots \times V_T$, where V_i is as in Definition E.4. By [64, Example 1.33], we have $\dim V \leq (d+1)T$. Since W is a variety of tensorization, by Lemma E.5, V_i is admissible with respect to $\{L_i(\cdot, \mathbf{W}) : \mathbf{W} \in W\}$.

We are now ready to apply Lemma E.1, which gives that when $T \geq \dim W$, there exists subvariety $Z \subset V$ with $\dim Z < \dim V \leq rT$, and for any $(\mathbf{X}_1, \dots, \mathbf{X}_T) \in V \setminus Z$ and any $\mathbf{W} \in W$, if $\langle \mathbf{X}_1, \mathbf{W} \rangle = \dots = \langle \mathbf{X}_T, \mathbf{W} \rangle = 0$, then $\mathbf{W} = 0$. By Lemma E.3, we have for every $(\mathbf{X}_1, \dots, \mathbf{X}_T) \in V \setminus Z$, the map $\mathbf{W} \mapsto (\langle \mathbf{X}_1, \mathbf{W} \rangle, \dots, \langle \mathbf{X}_T, \mathbf{W} \rangle)$ is injective, so \mathbf{W}_f and thus f can be uniquely recovered from the observed rewards.

Finally, we show that $(\varphi_1^{-1} \times \cdots \times \varphi_1^{-1})(Z)$ is a null set, where φ_1 is as in (10). According to the proof of Lemma E.1 by [69], we find that Z is also defined by homogeneous polynomials. We take the slice $Z' = \{\mathbf{x} \in Z : x_{11} = \cdots = x_{T1} = 1\}$, $V' = \{\mathbf{x} \in V : x_{11} = \cdots = x_{T1} = 1\}$, (here x_{ij} is the j -th coordinate of \mathbf{x}_i), then Z', V' are varieties. Since $\dim Z < \dim V$, we have $\dim Z' = \dim Z - T < \dim V - T = \dim V'$ and $Z' \subset V'$.

Now consider the regular map $\varphi'_1 : V' \rightarrow (\mathbb{C}^d)^T$,

$$\left(\binom{1}{\mathbf{x}_1}^{\otimes p}, \cdots, \binom{1}{\mathbf{x}_T}^{\otimes p} \right) \mapsto (\mathbf{x}_1, \cdots, \mathbf{x}_T). \quad (13)$$

Then $\varphi'_1(Z'), \varphi'_1(V')$ are both varieties. By [56, Lemma 9.9], we have $\dim \overline{\varphi'_1(Z')} \leq \dim Z$. Since $\varphi'_1(V) = (\mathbb{C}^d)^T$ and $\dim \varphi'_1(V') \leq \dim \overline{V'} = \dim V - T \leq (d+1)T - T$, we have $\dim V = \dim \varphi'_1(V) = dT$ and as a result, $\dim \overline{\varphi'_1(Z)} \leq \dim Z < \dim V = dT$. By Lemma E.6, $\overline{\varphi'_1(Z)}$ is a null set. Since $(\mathbf{x}_1, \cdots, \mathbf{x}_T) \notin \overline{\varphi'_1(Z)}$ implies that $(\varphi_1(\mathbf{x}_1), \cdots, \varphi_1(\mathbf{x}_T)) \notin Z$, we conclude the proof. \square

Theorem E.7 is stated for complex sample points. Next we extend it to the real case.

Lemma E.8. *In Lemma E.1, if we assume in addition that $\dim_{\mathbb{R}} V_{\mathbb{R}} = \dim V$, then the conclusion can be enhanced to ensure that Z is a real subvariety and $\dim_{\mathbb{R}} Z < \dim_{\mathbb{R}} V_{\mathbb{R}}$.*

Lemma E.9. *Let $V \subset \mathbb{R}^n$ be a (Zariski) closed proper subset, $V \neq \mathbb{R}^n$. Then V is a null set.*

The proof of Lemma E.9 is the same as that of Lemma E.6.

Theorem E.10. *We can additionally assume $\mathbf{x}_i \in \mathbb{R}^d$ in Theorem E.7.*

Proof. We verify that $\dim V = \dim_{\mathbb{R}} V_{\mathbb{R}}$, where V is defined in the proof of Theorem E.7, but this follows clearly by [12, Corollary 2.8.2]. We conclude the proof by applying Lemma E.8. \square

Finally, we apply Theorem E.10 to two concrete classes of polynomials, namely Examples 3.22 and 3.23. For Example 3.22, we construct its variety of tensorization of \mathcal{R}_{γ} as follows. We first construct the tensorization of each polynomial. We define

$$\mathbf{w}_f = \sum_{i=1}^r a_i \binom{1}{\mathbf{w}_i}^{\otimes p_i} \otimes \binom{1}{0}^{\otimes (p-p_i)}. \quad (14)$$

Next we construct the variety of tensorization W . Consider the map $\varphi_2 : (\mathbb{C}^d)^r \rightarrow \mathbb{C}^{(d+1)^p}$,

$$\varphi_2(\mathbf{w}_1, \cdots, \mathbf{w}_r) = \sum_{i=1}^r \binom{1}{\mathbf{w}_i}^{\otimes p_i} \otimes \binom{1}{0}^{\otimes (p-p_i)}, \quad (15)$$

and let $Y = \mathbb{P}(\overline{\text{Im } \varphi_2})$. Similar to V_i , we can prove that Y is an irreducible closed variety defined by homogeneous polynomials with $\dim Y \leq dr + 1$. Next consider the map $\varphi'_2 : (\mathbb{C}^d)^{2r} \rightarrow \mathbb{C}^{(d+1)^p}$,

$$\varphi'_2(\mathbf{w}_1, \cdots, \mathbf{w}_{2r}) = \varphi_2(\mathbf{w}_1, \cdots, \mathbf{w}_r) - \varphi_2(\mathbf{w}_{r+1}, \cdots, \mathbf{w}_{2r}) \quad (16)$$

and let $W = \mathbb{P}(\overline{\text{Im } \varphi'_2})$. Similar to Y , we can prove that W is an irreducible closed variety defined by homogeneous polynomials with $\dim W \leq 2dr + 1$. Together with Theorem E.10, we can conclude that the optimal action for Example 3.22 can be uniquely determined using at most $2dr + 1$ samples.

For Example 3.23, we construct W as follows. Let

$$\mathbf{U} = (\mathbf{w}_1 \cdots \mathbf{w}_k), \quad q = \sum_{I \subseteq [k]: |I| \leq p} a_I x^I,$$

then we construct the tensorization of each polynomial by

$$\mathbf{w}_f = \sum_{I \subseteq [k]: |I| \leq p} a_I \bigotimes_{i \in I} \binom{1}{\mathbf{w}_i} \otimes \binom{1}{0}^{\otimes (p-|I|)}. \quad (17)$$

Then we have $f(\mathbf{x}) = \langle \mathbf{W}_f, \mathbf{X}_x \rangle$. To reduce the dimension of W and get better sample complexity bound, we construct in a manner slightly different from what we did for Example 3.22. Consider the map $\varphi_3 : (\mathbb{C}^d)^k \times \mathbb{C}^{(k+1)^p} \rightarrow \mathbb{C}^{(d+1)^p}$,

$$(\mathbf{w}_1, \dots, \mathbf{w}_k) \times (a_I : I \subseteq [k], |I| \leq p) \mapsto \mathbf{W}_f, \quad (18)$$

where \mathbf{W}_f is as defined in (17). Let $Y = \mathbb{P}(\overline{\text{Im } \varphi_3})$ and $W = \mathbb{P}(\overline{\text{Im } \varphi_3} - \overline{\text{Im } \varphi_3})$. We end up with $\dim Y \leq dk + (k+1)^p + 1$, $\dim W \leq 2(dk + (k+1)^p) + 1$. So we conclude that the optimal action for Example 3.23 can be uniquely determined using at most $2dk + 2(k+1)^p + 1$ samples.

F Omitted Proof for Lower Bounds with UCB Algorithms

In this section, we provide the proof for the lower bounds for learning with UCB algorithms in Subsection G.0.1.

Notation Recall that we use Λ to denote the subset of the p -th multi-indices $\Lambda = \{(\alpha_1, \dots, \alpha_p) \mid 1 \leq \alpha_1 < \dots < \alpha_p \leq d\}$. For an $\alpha = (\alpha_1, \dots, \alpha_p) \in \Lambda$, denote $\mathbf{M}_\alpha = \mathbf{e}_{\alpha_1} \otimes \dots \otimes \mathbf{e}_{\alpha_p}$, $\mathbf{A}_\alpha = (\mathbf{e}_{\alpha_1} + \dots + \mathbf{e}_{\alpha_p})^{\otimes p}$. The model space \mathcal{M} is a subset of rank-1 p -th order tensors, which is defined as $\mathcal{M} = \{\mathbf{M}_\alpha \mid \alpha \in \Lambda\}$. We define the core action set \mathcal{A}_0 as $\mathcal{A}_0 = \{\mathbf{e}_{\alpha_1} + \dots + \mathbf{e}_{\alpha_p} \mid \alpha \in \Lambda\}$. The action set \mathcal{A} is the convex hull of \mathcal{A}_0 : $\mathcal{A} = \text{conv}(\mathcal{A}_0)$. Assume that the ground-truth parameter is $\mathbf{M}^* = \mathbf{M}_{\alpha^*} \in \mathcal{M}$. At round t , the algorithm chooses an action $\mathbf{a}_t \in \mathcal{A}$, and gets the **noiseless** reward $r_t = r(\mathbf{M}^*, \mathbf{a}_t) = \langle \mathbf{M}^*, (\mathbf{a}_t)^{\otimes p} \rangle = \prod_{i=1}^p \langle \mathbf{e}_{\alpha_i^*}, \mathbf{a}_t \rangle$.

F.1 Proof for Theorem G.2

We introduce a lemma showing that if the action set is **restricted** to the core action set \mathcal{A}_0 , then at least $|\mathcal{A}_0| - 1 = \binom{d}{p} - 1$ actions are needed to identify the ground-truth.

Lemma F.1. *If the actions are restricted to \mathcal{A}_0 , then for the noiseless degree- p polynomial bandits, any algorithm needs to play at least $\binom{d}{p} - 1$ actions to determine \mathbf{M}^* in the worst case. Furthermore, the worst-case cumulative regret at round T can be lower bounded by*

$$\mathfrak{R}(T) \geq \min\{T, \binom{d}{p} - 1\}.$$

proof of Lemma F.1. For any α and α' , the reward of playing $\mathbf{e}_{\alpha_1} + \dots + \mathbf{e}_{\alpha_p}$ when the ground-truth model is \mathbf{M}'_α is

$$\begin{aligned} \langle \mathbf{M}'_\alpha, (\mathbf{e}_{\alpha_1} + \dots + \mathbf{e}_{\alpha_p})^{\otimes p} \rangle &= \prod_{i=1}^p \langle \mathbf{e}_{\alpha'_i}, \mathbf{e}_{\alpha_1} + \dots + \mathbf{e}_{\alpha_p} \rangle \\ &= \prod_{i=1}^p \mathbb{I}\{\alpha'_i \in \alpha\} \\ &= \begin{cases} 1, & \text{if } \alpha = \alpha' \\ 0, & \text{otherwise.} \end{cases} \end{aligned}$$

Hence, no matter how the algorithm adaptively chooses the actions, in the worst case $\binom{d}{p} - 1$ actions are needed to determine \mathbf{M}^* . Also notice that the reward for $\mathbf{e}_{\alpha_1} + \dots + \mathbf{e}_{\alpha_p}$ is zero if $\alpha \neq \alpha^*$. Therefore the regret lower bound follows. \square

Next, we show that even when the action set is unrestricted, any UCB algorithm fails to explore in an unrestricted way. This is because the optimistic mechanism forbids the algorithm to play an informative action that is known to be low reward for all models in the confidence set. We first recall the definition of UCB algorithms.

UCB Algorithms The UCB algorithms sequentially maintain a confidence set \mathcal{C}_t after playing actions $\mathbf{a}_1, \dots, \mathbf{a}_t$. Then UCB algorithms play

$$\mathbf{a}_{t+1} \in \arg \max_{\mathbf{a} \in \mathcal{A}} \text{UCB}_t(\mathbf{a}),$$

where

$$\text{UCB}_t(\mathbf{a}) = \max_{\mathbf{M} \in \mathcal{C}_t} \langle \mathbf{M}, (\mathbf{a})^{\otimes p} \rangle.$$

proof of Theorem G.2. We prove that even if the action set is unrestricted, the optimistic mechanism in the UCB algorithm above forces it to choose actions in the restricted action set \mathcal{A}_0 .

Assume $\mathbf{M}^* = \mathbf{M}_{\alpha^*}$. Next we show that for all $\mathbf{a} \in \mathcal{A} - \mathcal{A}_0$ (where the minus sign should be understood as set difference), we have

$$\text{UCB}_t(\mathbf{a}) < 1.$$

For all $\mathbf{a} \in \mathcal{A}$, since $\mathcal{A} = \text{conv}(\mathcal{A}_0)$, we can write

$$\mathbf{a} = \sum_{\alpha \in \Lambda} p_\alpha (\mathbf{e}_{\alpha_1} + \dots + \mathbf{e}_{\alpha_p}),$$

where $\sum_{\alpha \in \Lambda} p_\alpha = 1$ and $p_\alpha \geq 0$. Therefore,

$$\begin{aligned} \text{UCB}_t(\mathbf{a}) &= \max_{\mathbf{M} \in \mathcal{C}_t} \langle \mathbf{M}, (\mathbf{a})^{\otimes p} \rangle \\ &\leq \max_{\mathbf{M} \in \mathcal{M}} \langle \mathbf{M}, (\mathbf{a})^{\otimes p} \rangle \\ &= \max_{\alpha'} \langle \mathbf{M}_{\alpha'}, (\mathbf{a})^{\otimes p} \rangle \\ &= \max_{\alpha'} \prod_{i=1}^p \langle \mathbf{e}_{\alpha'_i}, \mathbf{a} \rangle. \end{aligned}$$

Plug in the expression of \mathbf{a} , we have

$$\begin{aligned} \langle \mathbf{e}_{\alpha'_i}, \mathbf{a} \rangle &= \sum_{\alpha} p_\alpha \langle \mathbf{e}_{\alpha'_i}, \mathbf{e}_{\alpha_1} + \dots + \mathbf{e}_{\alpha_p} \rangle \\ &= \sum_{\alpha} p_\alpha \mathbb{I}\{\alpha'_i \in \alpha\} \\ &\leq \sum_{\alpha} p_\alpha = 1. \end{aligned}$$

Therefore, for any fixed $\alpha' = (\alpha'_1, \dots, \alpha'_p)$,

$$\begin{aligned} \prod_{i=1}^p \langle \mathbf{e}_{\alpha'_i}, \mathbf{a} \rangle &= \left(\sum_{\alpha} p_\alpha \mathbb{I}\{\alpha'_1 \in \alpha\} \right) \cdots \left(\sum_{\alpha} p_\alpha \mathbb{I}\{\alpha'_p \in \alpha\} \right) \\ &\leq 1, \end{aligned}$$

where the equality holds if and only if for any $p_\alpha > 0$, $\alpha = \alpha'$, which is equivalent to $\mathbf{a} = \mathbf{e}_{\alpha'_1} + \dots + \mathbf{e}_{\alpha'_p}$. Therefore, if $\mathbf{a} \in \mathcal{A} - \mathcal{A}_0$, for any $\alpha' \in \Lambda$, we have $\prod_{i=1}^p \langle \mathbf{e}_{\alpha'_i}, \mathbf{a} \rangle < 1$. This means

$$\text{UCB}_t(\mathbf{a}) < 1.$$

Meanwhile, we can see that for the action $\mathbf{a}^* = \mathbf{e}_{\alpha_1^*} + \dots + \mathbf{e}_{\alpha_p^*} \in \mathcal{A}_0$,

$$\begin{aligned} \text{UCB}_t(\mathbf{a}^*) &= \max_{\mathbf{M} \in \mathcal{C}_t} \langle \mathbf{M}, (\mathbf{a}^*)^{\otimes p} \rangle \\ &\geq \langle \mathbf{M}^*, (\mathbf{a}^*)^{\otimes p} \rangle && (\mathbf{M}^* \in \mathcal{C}_t) \\ &= \langle \mathbf{M}^*, \mathbf{A}_{\alpha^*} \rangle = 1. \end{aligned}$$

Therefore, we see that $(\mathcal{A} - \mathcal{A}_0) \cap \arg \max_{\mathbf{a} \in \mathcal{A}} \text{UCB}_t(\mathbf{a}) = \emptyset$, which means $\mathbf{a}_{t+1} \in \mathcal{A}_0$ for all $t \geq 0$. Therefore, by Lemma F.1, the theorem holds. \square

F.2 $O(d)$ Actions via Solving Polynomial Equations

Firstly, we verify that the model falls into the category of Example 3.23 with $k = p$. For every $\alpha \in \Lambda$, the reward of playing \mathbf{a} when the ground-truth model is \mathbf{M}_α is

$$\langle \mathbf{M}_\alpha, (\mathbf{a})^{\otimes p} \rangle = \prod_{i=1}^p \langle \mathbf{e}_{\alpha_i}, \mathbf{a} \rangle,$$

which can be written as $q_0(\mathbf{U}_\alpha \mathbf{a})$, where $q_0(x_1, \dots, x_p) = x_1 x_2 \cdots x_p$ and $\mathbf{U}_\alpha \in \mathbb{R}^{p \times d}$ is a matrix with \mathbf{e}_{α_i} as the i -th row.

Secondly, we show that since the ground-truth model is p -homogenous, we can extend the action set to $\text{conv}(\mathcal{A}, \mathbf{0})$. This is because for every action of the form $c\mathbf{a}$, where $0 \leq c \leq 1$ and $\mathbf{a} \in \mathcal{A}$, the reward is c^p times the reward at \mathbf{a} . Therefore, to get the reward at $c\mathbf{a}$, we only need to play at \mathbf{a} and multiply the reward by c^p .

Notice that $\text{conv}(\mathcal{A}, \mathbf{0})$ is of positive Lebesgue measure. By Theorem 3.21, we know that only $2(dk + (p+1)^p) = O(d)$ actions are needed to determine the optimal action almost surely.

G Proof of Section 3.3.2

We present the proof of Theorem 3.18 in the following.

Proof. We overload the notation and use $[d]$ to denote the set $\{e_1, e_2, \dots, e_d\}$. The hard instances are chosen in $\Delta \cdot [d]^p$, i.e. $(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p) = \Delta \cdot (\widehat{\boldsymbol{\theta}}_1, \dots, \widehat{\boldsymbol{\theta}}_p)$ where $(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p) \in [d]^p$. For a group of vectors $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p \in [d]$, we use

$$\text{supp}(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p) := (\max_{i \in [p]}(\boldsymbol{\theta}_i)_1, \dots, \max_{i \in [p]}(\boldsymbol{\theta}_i)_d) \in \{0, 1\}^d$$

to denote the support of these vectors. We use $\mathbf{a}^{(t)} \in \mathbb{R}^d$ to denote the action in t -th episode.

We use $\mathbb{P}_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p)}$ to denote the measure on outcomes induced by the interaction of the fixed policy and the bandit parameterised by $r = \prod_{i=1}^p (\boldsymbol{\theta}_i^\top \mathbf{a}) + \epsilon$. Specifically, We use \mathbb{P}_0 to denote the measure on outcomes induced by the interaction of the fixed policy and the pure noise bandit $r = \epsilon$.

$$\begin{aligned} & \mathfrak{R}(d, p, T) \\ & \geq \frac{1}{d^p} \sum_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p) \in \Delta \cdot [d]^p} \mathbb{E}_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p)} \left[T \Delta^p / p^{p/2} - \sum_{t=1}^T \prod_{i=1}^p (\boldsymbol{\theta}_i^\top \mathbf{a}^{(t)}) \right] \\ & = \frac{\Delta^p}{d^p} \sum_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p) \in \Delta \cdot [d]^p} \left(T / p^{p/2} - \mathbb{E}_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p)} \left[\sum_{t=1}^T \prod_{i=1}^p (\widehat{\boldsymbol{\theta}}_i^\top \mathbf{a}^{(t)}) \right] \right) \\ & \geq \frac{\Delta^p}{d^p} \sum_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p) \in \Delta \cdot [d]^p} \left(T / p^{p/2} - \mathbb{E}_0 \left[\sum_{t=1}^T \prod_{i=1}^p (\widehat{\boldsymbol{\theta}}_i^\top \mathbf{a}^{(t)}) \right] - T \|\mathbb{P}_0 - \mathbb{P}_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p)}\|_{\text{TV}} \right) \\ & \geq \frac{\Delta^p}{d^p} \sum_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p) \in \Delta \cdot [d]^p} \left(T / p^{p/2} - \mathbb{E}_0 \left[\sum_{t=1}^T \prod_{i=1}^p (\widehat{\boldsymbol{\theta}}_i^\top \mathbf{a}^{(t)}) \right] - T \sqrt{D_{\text{KL}}(\mathbb{P}_0 \| \mathbb{P}_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p)})} \right) \\ & \geq \frac{\Delta^p}{d^p} \sum_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p) \in \Delta \cdot [d]^p} \left(T / p^{p/2} - \mathbb{E}_0 \left[\sum_{t=1}^T \prod_{i=1}^p (\widehat{\boldsymbol{\theta}}_i^\top \mathbf{a}^{(t)}) \right] - T \sqrt{\Delta^{2p} \mathbb{E}_0 \left[\sum_{t=1}^T \prod_{i=1}^p (\widehat{\boldsymbol{\theta}}_i^\top \mathbf{a}^{(t)})^2 \right]} \right) \\ & \geq \frac{\Delta^p}{d^p} \left(\frac{d^p T}{p^{\frac{p}{2}}} - \mathbb{E}_0 \left[\sum_{(\widehat{\boldsymbol{\theta}}_1, \dots, \widehat{\boldsymbol{\theta}}_p) \in [d]^p} \sum_{t=1}^T \prod_{i=1}^p (\widehat{\boldsymbol{\theta}}_i^\top \mathbf{a}^{(t)}) \right] - T d^{\frac{p}{2}} \Delta^p \sqrt{\mathbb{E}_0 \left[\sum_{(\widehat{\boldsymbol{\theta}}_1, \dots, \widehat{\boldsymbol{\theta}}_p) \in [d]^p} \sum_{t=1}^T \prod_{i=1}^p (\widehat{\boldsymbol{\theta}}_i^\top \mathbf{a}^{(t)})^2 \right]} \right) \end{aligned}$$

where the first step comes from

$$\text{Regret} \geq \mathbb{E}_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p)} \left[T\Delta^p/p^{p/2} - \sum_{t=1}^T \prod_{i=1}^p (\boldsymbol{\theta}_i^\top \mathbf{a}^{(t)}) \right]$$

(the optimal action in hindsight is $\mathbf{a} = \text{supp}(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p)/\sqrt{p}$); the second step comes from $(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p) = \Delta \cdot (\widehat{\boldsymbol{\theta}}_1, \dots, \widehat{\boldsymbol{\theta}}_p)$ and algebra; the third step comes from $\left| \sum_{t=1}^T \prod_{i=1}^p (\widehat{\boldsymbol{\theta}}_i^\top \mathbf{a}^{(t)}) \right| \leq T$; the fourth step comes from Pinsker's inequality; the fifth step comes from

$$\begin{aligned} D_{\text{KL}}(\mathbb{P}_0 \parallel \mathbb{P}_{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p}) &= \mathbb{E}_0 \left[\sum_{t=1}^T D_{\text{KL}} \left(N(0, 1) \parallel N \left(\prod_{i=1}^p (\boldsymbol{\theta}_i^\top \mathbf{a}^{(t)}), 1 \right) \right) \right] \\ &= \Delta^{2p} \mathbb{E}_0 \left[\sum_{t=1}^T \prod_{i=1}^p (\widehat{\boldsymbol{\theta}}_i^\top \mathbf{a}^{(t)})^2 \right] \end{aligned}$$

and the final step comes from Jensen's inequality and algebra.

Notice that

$$\begin{aligned} \mathbb{E}_0 \left[\sum_{(\widehat{\boldsymbol{\theta}}_1, \dots, \widehat{\boldsymbol{\theta}}_p) \in [d]^p} \sum_{t=1}^T \prod_{i=1}^p (\widehat{\boldsymbol{\theta}}_i^\top \mathbf{a}^{(t)}) \right] &= \mathbb{E}_0 \left[\sum_{(j_1, \dots, j_p) \in [d]^p} \sum_{t=1}^T \prod_{i=1}^p (\mathbf{a}_{j_i}^{(t)}) \right] \\ &= \mathbb{E}_0 \left[\sum_{t=1}^T \prod_{i=1}^p \left(\sum_{j=1}^d \mathbf{a}_j^{(t)} \right) \right] \\ &\leq \mathbb{E}_0 \left[\sum_{t=1}^T \prod_{i=1}^p \|\mathbf{a}^{(t)}\|_1 \right] \\ &\leq d^{p/2} T \end{aligned}$$

and

$$\begin{aligned} \mathbb{E}_0 \left[\sum_{(\widehat{\boldsymbol{\theta}}_1, \dots, \widehat{\boldsymbol{\theta}}_p) \in [d]^p} \sum_{t=1}^T \prod_{i=1}^p (\widehat{\boldsymbol{\theta}}_i^\top \mathbf{a}^{(t)})^2 \right] &= \mathbb{E}_0 \left[\sum_{(j_1, \dots, j_p) \in [d]^p} \sum_{t=1}^T \prod_{i=1}^p (\mathbf{a}_{j_i}^{(t)})^2 \right] \\ &= \mathbb{E}_0 \left[\sum_{t=1}^T \prod_{i=1}^p \|\mathbf{a}^{(t)}\|_2^2 \right] \\ &\leq T \end{aligned}$$

where we used $\|\mathbf{a}^{(t)}\|_2 \leq 1, \forall t \in [T]$. Therefore plugging back we have

$$\mathfrak{R}(d, p, T) \geq \frac{\Delta^p}{d^p} \left(\frac{d^p T}{p^{p/2}} - d^{p/2} T - T d^{p/2} \Delta^p \sqrt{T} \right)$$

and finally letting $\Delta^p = \sqrt{\frac{d^p}{4Tp^p}}$ leads to

$$\mathfrak{R}(d, p, T) \geq O(\sqrt{d^p T}/p^p).$$

□

Remark G.1. Better result $O(\sqrt{d^p T})$ holds for bandits $r = \prod_{i=1}^p (\boldsymbol{\theta}_i^\top \mathbf{a}_i) + \epsilon$ where $\mathbf{a}_i \in \mathbb{R}^d, \|\mathbf{a}_i\|_2 \leq 1$.

For completeness, we show the proof of the above remark.

Proof. We overload the notation and use $[d]$ to denote the set $\{e_1, e_2, \dots, e_d\}$. The hard instances are chosen in $\Delta \cdot [d]^p$, i.e. $(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p) = \Delta \cdot (\widehat{\boldsymbol{\theta}}_1, \dots, \widehat{\boldsymbol{\theta}}_p)$ where $(\widehat{\boldsymbol{\theta}}_1, \dots, \widehat{\boldsymbol{\theta}}_p) \in [d]^p$. We use $\mathbf{a}_i^{(t)} \in \mathbb{R}^d$ to denote the i -th action in t -th episode, where $i \in [p], t \in [T]$.

We use $\mathbb{P}_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p)}$ to indicate the measure on outcomes induced by the interaction of the fixed policy and the bandit parameterised by $r = \prod_{i=1}^p (\boldsymbol{\theta}_i^\top \mathbf{a}_i) + \epsilon$. Specifically, We use \mathbb{P}_0 to indicate the measure on outcomes induced by the interaction of the fixed policy and the pure noise bandit $r = \epsilon$.

$$\begin{aligned}
& \mathfrak{R}(d, p, T) \\
& \geq \frac{1}{d^p} \sum_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p) \in \Delta \cdot [d]^p} \mathbb{E}_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p)} \left[T\Delta^p - \sum_{t=1}^T \prod_{i=1}^p (\boldsymbol{\theta}_i^\top \mathbf{a}_i^{(t)}) \right] \\
& = \frac{\Delta^p}{d^p} \sum_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p) \in \Delta \cdot [d]^p} \left(T - \mathbb{E}_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p)} \left[\sum_{t=1}^T \prod_{i=1}^p (\widehat{\boldsymbol{\theta}}_i^\top \mathbf{a}_i^{(t)}) \right] \right) \\
& \geq \frac{\Delta^p}{d^p} \sum_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p) \in \Delta \cdot [d]^p} \left(T - \mathbb{E}_0 \left[\sum_{t=1}^T \prod_{i=1}^p (\widehat{\boldsymbol{\theta}}_i^\top \mathbf{a}_i^{(t)}) \right] - T \|\mathbb{P}_0 - \mathbb{P}_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p)}\|_{\text{TV}} \right) \\
& \geq \frac{\Delta^p}{d^p} \sum_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p) \in \Delta \cdot [d]^p} \left(T - \mathbb{E}_0 \left[\sum_{t=1}^T \prod_{i=1}^p (\widehat{\boldsymbol{\theta}}_i^\top \mathbf{a}_i^{(t)}) \right] - T \sqrt{D_{\text{KL}}(\mathbb{P}_0 \| \mathbb{P}_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p)})} \right) \\
& \geq \frac{\Delta^p}{d^p} \sum_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p) \in \Delta \cdot [d]^p} \left(T - \mathbb{E}_0 \left[\sum_{t=1}^T \prod_{i=1}^p (\widehat{\boldsymbol{\theta}}_i^\top \mathbf{a}_i^{(t)}) \right] - T \sqrt{\Delta^{2p} \mathbb{E}_0 \left[\sum_{t=1}^T \prod_{i=1}^p (\widehat{\boldsymbol{\theta}}_i^\top \mathbf{a}_i^{(t)})^2 \right]} \right) \\
& \geq \frac{\Delta^p}{d^p} \left(d^p T - \mathbb{E}_0 \left[\sum_{(\widehat{\boldsymbol{\theta}}_1, \dots, \widehat{\boldsymbol{\theta}}_p) \in [d]^p} \sum_{t=1}^T \prod_{i=1}^p (\widehat{\boldsymbol{\theta}}_i^\top \mathbf{a}_i^{(t)}) \right] - T d^{\frac{p}{2}} \Delta^p \sqrt{\mathbb{E}_0 \left[\sum_{(\widehat{\boldsymbol{\theta}}_1, \dots, \widehat{\boldsymbol{\theta}}_p) \in [d]^p} \sum_{t=1}^T \prod_{i=1}^p (\widehat{\boldsymbol{\theta}}_i^\top \mathbf{a}_i^{(t)})^2 \right]} \right)
\end{aligned}$$

where the first step comes from

$$\text{Regret} \geq \mathbb{E}_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p)} \left[T\Delta^p - \sum_{t=1}^T \prod_{i=1}^p (\boldsymbol{\theta}_i^\top \mathbf{a}_i^{(t)}) \right]$$

(the optimal action in hindsight is $\mathbf{a}_i = \widehat{\boldsymbol{\theta}}_i$); the second step comes from $(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p) = \Delta \cdot (\widehat{\boldsymbol{\theta}}_1, \dots, \widehat{\boldsymbol{\theta}}_p)$ and algebra; the third step comes from $\left| \sum_{t=1}^T \prod_{i=1}^p (\widehat{\boldsymbol{\theta}}_i^\top \mathbf{a}_i^{(t)}) \right| \leq T$; the fourth step comes from Pinsker's inequality; the fifth step comes from

$$\begin{aligned}
D_{\text{KL}}(\mathbb{P}_0 \| \mathbb{P}_{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p}) &= \mathbb{E}_0 \left[\sum_{t=1}^T D_{\text{KL}} \left(N(0, 1) \| N\left(\prod_{i=1}^p (\boldsymbol{\theta}_i^\top \mathbf{a}_i^{(t)}), 1\right) \right) \right] \\
&= \Delta^{2p} \mathbb{E}_0 \left[\sum_{t=1}^T \prod_{i=1}^p (\widehat{\boldsymbol{\theta}}_i^\top \mathbf{a}_i^{(t)})^2 \right]
\end{aligned}$$

and the final step comes from Jensen's inequality and algebra.

Notice that

$$\begin{aligned}
\mathbb{E}_0 \left[\sum_{(\widehat{\boldsymbol{\theta}}_1, \dots, \widehat{\boldsymbol{\theta}}_p) \in [d]^p} \sum_{t=1}^T \prod_{i=1}^p (\widehat{\boldsymbol{\theta}}_i^\top \mathbf{a}_i^{(t)}) \right] &= \mathbb{E}_0 \left[\sum_{(j_1, \dots, j_p) \in [d]^p} \sum_{t=1}^T \prod_{i=1}^p ((\mathbf{a}_i^{(t)})_{j_i}) \right] \\
&= \mathbb{E}_0 \left[\sum_{t=1}^T \prod_{i=1}^p \left(\sum_{j=1}^d (\mathbf{a}_i^{(t)})_j \right) \right] \\
&\leq \mathbb{E}_0 \left[\sum_{t=1}^T \prod_{i=1}^p \|\mathbf{a}_i^{(t)}\|_1 \right] \\
&\leq d^{p/2} T
\end{aligned}$$

and

$$\begin{aligned} \mathbb{E}_0 \left[\sum_{(\hat{\boldsymbol{\theta}}_1, \dots, \hat{\boldsymbol{\theta}}_p) \in [d]^p} \sum_{t=1}^T \prod_{i=1}^p (\hat{\boldsymbol{\theta}}_i^\top \mathbf{a}_i^{(t)})^2 \right] &= \mathbb{E}_0 \left[\sum_{(j_1, \dots, j_p) \in [d]^p} \sum_{t=1}^T \prod_{i=1}^p ((\mathbf{a}_i^{(t)})_{j_i})^2 \right] \\ &= \mathbb{E}_0 \left[\sum_{t=1}^T \prod_{i=1}^p \|\mathbf{a}_i^{(t)}\|_2^2 \right] \\ &\leq T \end{aligned}$$

where we used $\|\mathbf{a}_i^{(t)}\|_2 \leq 1, \forall t \in [T]$. Therefore plugging back we have

$$\mathfrak{R}(d, p, T) \geq \frac{\Delta^p}{d^p} \left(d^p T - d^{\frac{p}{2}} T - T d^{\frac{p}{2}} \Delta^p \sqrt{T} \right)$$

and finally letting $\Delta^p = \sqrt{\frac{d^p}{4T}}$ leads to

$$\mathfrak{R}(d, p, T) \geq O(\sqrt{d^p T}).$$

□

We present the proof of Theorem 3.19 in the following.

Proof. Denote the optimal action in hindsight as $\mathbf{a}^* = \text{supp}(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p) / \sqrt{p}$. From the proof of Theorem 3.18 we know that if $T \leq \frac{1}{4p^{\frac{p}{2}}} \cdot \frac{d^p}{\Delta^{2p}}$, then

$$\begin{aligned} &\frac{1}{d^p} \sum_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p) \in \Delta \cdot [d]^p} \mathbb{E}_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p)} \left[\prod_{i=1}^p (\boldsymbol{\theta}_i^\top \mathbf{a}^*) - \prod_{i=1}^p (\boldsymbol{\theta}_i^\top \mathbf{a}^{(t)}) \right] \\ &\geq \frac{\Delta^p}{d^p} \left(\frac{d^p}{p^{\frac{p}{2}}} - d^{\frac{p}{2}} - d^{\frac{p}{2}} \Delta^p \sqrt{T} \right) \\ &\geq \frac{\Delta^p}{4p^{\frac{p}{2}}} \\ &\geq \frac{1}{4} \cdot \frac{1}{d^p} \sum_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p) \in \Delta \cdot [d]^p} \mathbb{E}_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p)} \left[\prod_{i=1}^p (\boldsymbol{\theta}_i^\top \mathbf{a}^*) \right] \end{aligned}$$

which indicates the following

$$\inf_{\pi} \sup_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p)} \mathbb{E}_{(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p)} \left[\frac{3}{4} \cdot \prod_{i=1}^p (\boldsymbol{\theta}_i^\top \mathbf{a}^*) - \prod_{i=1}^p (\boldsymbol{\theta}_i^\top \mathbf{a}^{(t)}) \right] \geq 0.$$

□

G.0.1 Lower Bounds with UCB Algorithms

In this subsection, we construct a hard bandit problem where the rewards are noiseless degree- p polynomial, and show that any UCB algorithm needs at least $\Omega(d^p)$ actions to learn the optimal action. On the contrary, Theorem 3.21 shows that by playing actions randomly, we only need $2(dk + (p+1)^p) = O(d)$ actions.

Hard Case Construction Let \mathbf{e}_i denotes the i -th standard orthonormal basis of \mathbb{R}^d , i.e., \mathbf{e}_i has only one 1 at the i -th entry and 0's for other entries. We define a p -th multi-indices set $\Lambda = \{(\alpha_1, \dots, \alpha_p) | 1 \leq \alpha_1 < \dots < \alpha_p \leq d\}$. For an $\alpha = (\alpha_1, \dots, \alpha_p) \in \Lambda$, denote $\mathbf{M}_\alpha = \mathbf{e}_{\alpha_1} \otimes \dots \otimes \mathbf{e}_{\alpha_p}$. Then the model space \mathcal{M} is defined as $\mathcal{M} = \{\mathbf{M}_\alpha | \alpha \in \Lambda\}$, which is a subset of rank-1 p -th order tensors. The action set \mathcal{A} is defined as $\mathcal{A} = \text{conv}(\{\mathbf{e}_{\alpha_1} + \dots + \mathbf{e}_{\alpha_p} | \alpha \in \Lambda\})$. Assume that the ground-truth parameter is $\mathbf{M}^* = \mathbf{M}_{\alpha^*} \in \mathcal{M}$. The noiseless reward $r_t = r(\mathbf{M}^*, \mathbf{a}_t) = \langle \mathbf{M}^*, (\mathbf{a}_t)^{\otimes p} \rangle = \prod_{i=1}^p \langle \mathbf{e}_{\alpha_i^*}, \mathbf{a}_t \rangle$ is a polynomial of \mathbf{a}_t and falls into the case of Example 3.23.

UCB Algorithms The UCB algorithms sequentially maintain a confidence set \mathcal{C}_t after playing actions $\mathbf{a}_1, \dots, \mathbf{a}_t$. Then UCB algorithms play $\mathbf{a}_{t+1} \in \arg \max_{\mathbf{a} \in \mathcal{A}} \text{UCB}_t(\mathbf{a})$, where $\text{UCB}_t(\mathbf{a}) = \max_{\mathbf{M} \in \mathcal{C}_t} \langle \mathbf{M}, (\mathbf{a})^{\otimes p} \rangle$.

Theorem G.2. *Assume that for each $t \geq 0$, the confidence set \mathcal{C}_t contains the ground-truth model, i.e., $\mathbf{M}^* \in \mathcal{C}_t$. Then for the noiseless degree- p polynomial bandits, any UCB algorithm needs to play at least $\binom{d}{p} - 1$ actions to distinguish models in \mathcal{M} . Furthermore, the worst-case cumulative regret at round T can be lower bounded by*

$$\mathfrak{R}(T) \geq \min\left\{T, \binom{d}{p} - 1\right\}.$$

Theorem G.2 shows the failure of the optimistic mechanism, which forbids the algorithm to play an informative action that is known to be of low reward for all models in the confidence set. On the contrary, the reward function class falls into the form of $q(\mathbf{U}\mathbf{a})$, therefore, by playing actions randomly⁶, we only need $O(d)$ actions as Theorem 3.21 suggests.

⁶Careful readers may notice that \mathcal{A} is of measure zero in this setting. However, since the reward function is a homogenous polynomial of degree p , we can actually obtain the rewards on $\text{conv}(\mathcal{A}, \mathbf{0})$, which is of positive measure.