
Efficient Linear Bandits through Matrix Sketching

Ilja Kuzborskij

Leonardo Cella

Nicolò Cesa-Bianchi

Dipartimento di Informatica & DSRC, Università degli Studi di Milano, 20133 Milano, Italy

Abstract

We prove that two popular linear contextual bandit algorithms, OFUL and Thompson Sampling, can be made efficient using Frequent Directions, a deterministic online sketching technique. More precisely, we show that a sketch of size m allows a $\mathcal{O}(md)$ update time for both algorithms, as opposed to $\Omega(d^2)$ required by their non-sketched versions in general (where d is the dimension of context vectors). This computational speedup is accompanied by regret bounds of order $(1 + \varepsilon_m)^{3/2}d\sqrt{T}$ for OFUL and of order $((1 + \varepsilon_m)d)^{3/2}\sqrt{T}$ for Thompson Sampling, where ε_m is bounded by the sum of the tail eigenvalues not covered by the sketch. In particular, when the selected contexts span a subspace of dimension at most m , our algorithms have a regret bound matching that of their slower, non-sketched counterparts. Experiments on real-world datasets corroborate our theoretical results.

1 Introduction

The stochastic contextual bandit is a sequential decision-making problem where an agent interacts with an unknown environment in a series of rounds. In each round, the environment reveals a set of feature vectors (called contexts, or actions) to the agent. The agent chooses an action from the revealed set and observes the stochastic reward associated with that action (bandit feedback). The strategy used by the agent for choosing actions based on past observations is called a policy. The goal of the agent is to learn a policy minimizing the regret, defined as the difference between the total reward of the optimal policy (i.e., the policy choosing

the action with highest expected reward at each round) and the total reward of the agent’s policy.

Contextual bandits are a popular modelling tool in many interactive machine learning tasks. A typical area of application is personalized recommendation, where a recommender system selects a product for a given user from a set of available products (each described by a feature vector) and receives a feedback (purchase or non-purchase) for the selected product.

We focus on the *stochastic linear bandit* model (Auer, 2002; Dani et al., 2008), where the set of actions (or decision set) is a finite¹ set $D_t \subset \mathbb{R}^d$, and the reward for choosing action $\mathbf{x}_t \in D_t$ is given by $Y_t = \mathbf{x}_t^\top \mathbf{w}^* + \eta_t$ where $\mathbf{w}^* \in \mathbb{R}^d$ is a fixed and unknown vector of real coefficients and η_t is a zero-mean random variable. The regret in this setting is defined by

$$R_T = \sum_{t=1}^T \mathbf{x}_t^{*\top} \mathbf{w}^* - \sum_{t=1}^T \mathbf{x}_t^\top \mathbf{w}^* \quad (1)$$

where $\mathbf{x}_t^* = \arg \max_{\mathbf{x} \in D_t} \mathbf{x}^\top \mathbf{w}^*$ is the optimal action at round t . Bounds on the regret typically apply to any individual sequence of decision sets D_t and depend on quantities arising from the interplay between \mathbf{w}^* , the sequence of decision sets, and the randomness of the rewards. Note that R_T is a random variable because the actions $\mathbf{x}_t \in D_t$ selected by the policy are functions of the past observed rewards. For this reason, our regret bounds only hold with probability at least $1 - \delta$, where δ is a confidence parameter. By choosing $\delta = T^{-1}$, we can instead bound the expected regret $\mathbb{E}[R_T]$ by paying only a $\ln T$ extra factor in the bound.

We consider two of the most popular algorithms for stochastic linear bandits: OFUL (Abbasi-Yadkori et al., 2011) and linear Thompson Sampling (Agrawal and Goyal, 2013) (linear TS for short). While exhibiting good theoretical and empirical performances, both algorithms require $\Omega(d^2)$ time to update their model after each round. In this work we investigate whether it is possible to significantly reduce this update time while ensuring that the regret remains nicely bounded.

¹Note that our regret bounds do not actually depend on the cardinality of the sets D_t .

The quadratic dependence on d is due to the computation of the inverse correlation matrix of past actions (a cubic dependence is avoided because each new inverse is a rank-one perturbation of the previous inverse). The occurrence of this matrix is caused by the linear nature of rewards: to compute their decisions, both algorithms essentially solve a regularized least squares problem at every round. In order to improve the running time, we sketch the correlation matrix using a specific technique —Frequent Directions, (Ghashami et al., 2016)— that works well in a sequential learning setting. While matrix sketching is a well-known approach (Woodruff, 2014), to the best of our knowledge this is the first work that applies sketching to linear contextual bandits while providing rigorous performance guarantees.

With a sketch size of m , a rank-one update of the correlation matrix takes only time $\mathcal{O}(md)$, which is linear in d for a constant sketch size. However, this speed-up comes at a price, as sketching reduces the matrix rank causing a loss of information which —in turn— affects the least squares estimates used by the algorithms. Our main technical contribution shows that when OFUL and linear TS are run with a sketched correlation matrix, their regret blows up by a factor which is controlled by the spectral decay of the correlation matrix of selected actions. More precisely, we show that the sketched variant of OFUL, called SOFUL, achieves a regret bounded by

$$R_T \stackrel{\tilde{\mathcal{O}}}{=} (1 + \varepsilon_m)^{\frac{3}{2}} \left(m + d \ln(1 + \varepsilon_m) \right) \sqrt{T} \quad (2)$$

where m is the sketch size and ε_m is upper bounded by the spectral tail (sum of the last $d - m + 1$ eigenvalues) of the correlation matrix for all T rounds. In the special case when the selected actions span a number of dimensions equal or smaller than the sketch size, then $\varepsilon_m = 0$ implying a regret of order $m\sqrt{T}$. Thus, we have a regret bound matching that of the slower, non-sketched counterpart.² When the correlation matrix has rank larger than the sketch size, the regret of SOFUL remains small to the extent the spectral tail of the matrix grows slowly with T . In the worst case of a spectrum with heavy tails, SOFUL may incur linear regret. In this respect, sketching is only justified when the computational cost of running OFUL cannot be afforded. Similarly, we prove that the efficient sketched formulation of linear TS enjoys a regret bound of order

$$R_T \stackrel{\tilde{\mathcal{O}}}{=} \left(m + d \ln(1 + \varepsilon_m) \right) (1 + \varepsilon_m)^{\frac{3}{2}} \sqrt{dT} . \quad (3)$$

Once again, for $\varepsilon_m = 0$ our bound is of order $m\sqrt{dT}$, which matches the regret bound for linear TS. When

²The regret bound of OFUL in (Abbasi-Yadkori et al., 2011, Theorem 3) is stated as $\mathcal{O}(d\sqrt{T})$, however, it can be improved for low-rank problems by using the “log-det” formulation of the confidence ellipsoid.

the rank of the correlation matrix is larger than the sketch size, the bound for linear TS behaves similarly to the bound for SOFUL.

Finally, we show a problem-dependent regret bound for SOFUL. This bound, which exhibits a logarithmic dependence on T , depends on the smallest gap Δ between the expected reward of the best and the second best action across the T rounds,

$$R_T \stackrel{\tilde{\mathcal{O}}}{=} \frac{1}{\Delta} (1 + \varepsilon_m)^3 \left(m + d \ln(1 + \varepsilon_m) \right)^2 (\ln T)^2 . \quad (4)$$

When $\varepsilon_m(T) = 0$ this bound is of order $\frac{m^2}{\Delta} (\ln T)^2$ which matches the corresponding bound for OFUL. Experiments on six real-world datasets support our theoretical results.

Additional related work. For an introduction to contextual bandits, we refer the reader to the recent monograph of Lattimore and Szepesvári (2018). The idea of applying sketching techniques to linear contextual bandits was also investigated by Yu et al. (2017), where they used random projections to preliminarily draw a random m -dimensional subspace which is then used in every round of play. However, the per-step computation time of their algorithm is cubic in m rather than quadratic like ours. Moreover, random projection introduces an additive error ε in the instantaneous regret which becomes of order $m^{-1/2}$ for any value of the confidence parameter δ bounded away from 1. A different notion of compression in contextual bandits is explored by Jun et al. (2017), where they use hashing algorithms to obtain a computation time sublinear in the number K of actions. An application of sketching (including Frequent Directions) to speed up 2nd order algorithms for online learning is studied by Luo et al. (2016), in a RKHS setting by Calandriello et al. (2017), and in stochastic optimization by Gonen et al. (2016).

2 Notation and preliminaries

Let $\mathcal{B}(\mathbf{z}, r) \subset \mathbb{R}^d$ be the Euclidean ball of center \mathbf{z} and radius $r > 0$ and let $\mathcal{B}(r) = \mathcal{B}(\mathbf{0}, r)$. Given a positive definite $d \times d$ matrix \mathbf{A} , we define the inner product $\langle \mathbf{x}, \mathbf{z} \rangle_{\mathbf{A}} = \mathbf{x}^\top \mathbf{A} \mathbf{z}$ and the induced norm $\|\mathbf{x}\|_{\mathbf{A}} = \sqrt{\mathbf{x}^\top \mathbf{A} \mathbf{x}}$, for any $\mathbf{x}, \mathbf{z} \in \mathbb{R}^d$. Throughout the paper, we write $f \stackrel{\tilde{\mathcal{O}}}{=} g$ to denote $f = \tilde{\mathcal{O}}(g)$. The contextual bandit protocol is described in Algorithm 1.

Algorithm 1 (Contextual Bandit)

- 1: **for** $t = 1, 2, \dots$ **do**
 - 2: Get decision set $D_t \subset \mathbb{R}^d$
 - 3: Use current policy to select action $\mathbf{x}_t \in D_t$
 - 4: Observe reward $Y_t \in \mathbb{R}$
 - 5: Use pair (\mathbf{x}_t, Y_t) to update the current policy
 - 6: **end for**
-

We introduce some standard assumptions for the linear contextual bandit setting. At any round $t = 1, 2, \dots$ the decision set $D_t \subset \mathbb{R}^d$ is finite and such that $\|\mathbf{x}\| \leq L$ for all $\mathbf{x} \in D_t$ and for all $t \geq 1$. The noise sequence $\eta_1, \eta_2, \dots, \eta_T$ is conditionally R -subgaussian for some fixed constant $R \geq 0$. Formally, for all $t \geq 1$ and all $\lambda \in \mathbb{R}$, $\mathbb{E}[e^{\lambda \eta_t} \mid \eta_1, \dots, \eta_{t-1}] \leq \exp(\lambda^2 R^2 / 2)$. Note that this implies $\mathbb{E}[\eta_t \mid \eta_1, \dots, \eta_{t-1}] = 0$ and $\text{Var}[\eta_t \mid \eta_1, \dots, \eta_{t-1}] \leq R^2$. Finally, we assume that a known upper bound S on $\|\mathbf{w}^*\|$ is available.

Both OFUL and Linear TS operate by computing a confidence ellipsoid to which \mathbf{w}^* belongs with high probability. Let $\mathbf{X}_t = [\mathbf{x}_1, \dots, \mathbf{x}_t]^\top$ be the $t \times d$ matrix of all actions selected up to round t by an arbitrary policy for linear contextual bandits. For $\lambda > 0$, define the regularized correlation matrix of actions \mathbf{V}_t and the regularized least squares (RLS) estimate $\hat{\mathbf{w}}_t$ as

$$\mathbf{V}_t = \mathbf{X}_t^\top \mathbf{X}_t + \lambda \mathbf{I} \quad \text{and} \quad \hat{\mathbf{w}}_t = \mathbf{V}_t^{-1} \sum_{s=1}^t \mathbf{x}_s Y_s. \quad (5)$$

The following theorem (Abbasi-Yadkori et al., 2011, Theorem 2) bounds in probability the distance, in terms of the norm $\|\cdot\|_{\mathbf{V}_t}$, between the optimal parameter \mathbf{w}^* and the RLS estimate $\hat{\mathbf{w}}_t$.

Theorem 1 (Confidence Ellipsoid). *Let $\hat{\mathbf{w}}_t$ be the RLS estimate constructed by an arbitrary policy for linear contextual bandits after t rounds of play. For any $\delta \in (0, 1)$, the optimal parameter \mathbf{w}^* belongs to the set $C_t \equiv \{\mathbf{w} \in \mathbb{R}^d : \|\mathbf{w} - \hat{\mathbf{w}}_t\|_{\mathbf{V}_t} \leq \beta_t(\delta)\}$ with probability at least $1 - \delta$, where*

$$\beta_t(\delta) = R \sqrt{d \ln \left(1 + \frac{tL^2}{\lambda d} \right) + 2 \ln \left(\frac{1}{\delta} \right) + S \sqrt{\lambda}}. \quad (6)$$

OFUL. The actions selected by OFUL are solutions to the following constrained optimization problem

$$\begin{aligned} \mathbf{x}_t &= \arg \max_{\mathbf{x} \in D_t} \max_{\mathbf{w} \in \mathbb{R}^d} \mathbf{x}^\top \mathbf{w} \\ \text{such that} \quad & \|\mathbf{w} - \hat{\mathbf{w}}_{t-1}\|_{\mathbf{V}_{t-1}} \leq \beta_{t-1}(\delta). \end{aligned}$$

Using Lemma 5 (in the appendix), OFUL can be formulated as Algorithm 2. Note that \mathbf{x}_t maximizes the expected reward estimate $\hat{\mathbf{w}}_{t-1}^\top \mathbf{x}$ plus a term $\beta_{t-1}(\delta) \|\mathbf{x}\|_{\mathbf{V}_{t-1}^{-1}}$ that provides an upper confidence bound for the RLS estimate in the direction of \mathbf{x} .

Linear TS. The linear Thompson Sampling algorithm of Agrawal and Goyal (2013) is Bayesian in nature: the selected actions and the observed rewards are used to update a Gaussian prior over the parameter space. Each action \mathbf{x}_t is selected by maximizing $\mathbf{x}^\top \hat{\mathbf{w}}_t^{\text{TS}}$ over $\mathbf{x} \in D_t$, where $\hat{\mathbf{w}}_t^{\text{TS}}$ is a random vector drawn from the posterior. As shown by Abeille and Lazaric (2017),

Algorithm 2 (OFUL)

Input: $\delta, \lambda > 0$

- 1: $\hat{\mathbf{w}}_0 = \mathbf{0}, \mathbf{V}_0^{-1} = \frac{1}{\lambda} \mathbf{I}$.
 - 2: **for** $t = 1, 2, \dots$ **do**
 - 3: Get decision set D_t
 - 4: Play $\mathbf{x}_t \leftarrow \arg \max_{\mathbf{x} \in D_t} \left\{ \hat{\mathbf{w}}_{t-1}^\top \mathbf{x} + \beta_{t-1}(\delta) \|\mathbf{x}\|_{\mathbf{V}_{t-1}^{-1}} \right\}$
 - 5: Observe reward Y_t
 - 6: Compute \mathbf{V}_t^{-1} and $\hat{\mathbf{w}}_t$ using (5)
 - 7: **end for**
-

linear TS can be equivalently defined as a randomized algorithm based on the RLS estimate (see Algorithm 3). The random vectors \mathbf{Z}_t are drawn i.i.d. from a suitable

Algorithm 3 (Linear TS)

Input: $\delta, \lambda > 0, m \in \{1, \dots, d-1\}$, \mathcal{D}^{TS} (sampling distribution)

- 1: $\hat{\mathbf{w}}_0 = \mathbf{0}, \mathbf{V}_0^{-1} = \frac{1}{\lambda} \mathbf{I}_{d \times d}, \delta' = \delta / (4T)$
 - 2: **for** $t = 1, 2, \dots$ **do**
 - 3: Get decision set D_t
 - 4: Sample $\mathbf{Z}_t \sim \mathcal{D}^{\text{TS}}$
 - 5: Play $\mathbf{x}_t \leftarrow \arg \max_{\mathbf{x} \in D_t} \mathbf{x}^\top \left(\hat{\mathbf{w}}_{t-1} + \tilde{\beta}_t(\delta') \mathbf{V}_{t-1}^{-\frac{1}{2}} \mathbf{Z}_t \right)$
 - 6: Observe reward Y_t
 - 7: Compute $\mathbf{V}_t^{-\frac{1}{2}}$ and $\hat{\mathbf{w}}_t$ using (5)
 - 8: **end for**
-

multivariate distribution \mathcal{D}^{TS} that need not be related to the posterior. In order to prove regret bounds, it is sufficient that the law of \mathbf{Z}_t satisfies certain properties.

Definition 1 (TS-sampling distribution). *A multivariate distribution \mathcal{D}^{TS} on \mathbb{R}^d , absolutely continuous w.r.t. the Lebesgue measure, is TS-sampling if it satisfies the following two properties:*

- (Anti-concentration) *There exists $p > 0$ such that for any \mathbf{u} with $\|\mathbf{u}\| = 1$, $\mathbb{P}(\mathbf{u}^\top \mathbf{Z} \geq 1) \geq p$.*
- (Concentration) *There exist $c, c' > 0$ such that for all $\delta \in (0, 1)$,*

$$\mathbb{P} \left(\|\mathbf{Z}\| \leq \sqrt{cd \ln \left(\frac{c'd}{\delta} \right)} \right) \geq 1 - \delta.$$

Similarly to OFUL, linear TS uses the notion of confidence ellipsoid. However, due to the properties of the sampling distribution \mathcal{D}^{TS} , the ellipsoid used by linear TS is larger by a factor of order \sqrt{d} than the ellipsoid used by OFUL. This causes an extra factor of \sqrt{d} in the regret bound, which is not known to be necessary.

Note that both OFUL and linear TS need to maintain \mathbf{V}_t^{-1} (or $\mathbf{V}_t^{-\frac{1}{2}}$), which requires time $\Omega(d^2)$ to update. In the next section, we show how this update time can be improved by sketching the regularized correlation matrix \mathbf{V}_t .

3 Sketching the correlation matrix

The idea of sketching is to maintain an approximation of \mathbf{X}_t , denoted by $\mathbf{S}_t \in \mathbb{R}^{m \times d}$, where $m \ll d$ is a small constant called the sketch size. If we choose m such that $\mathbf{S}_t^\top \mathbf{S}_t$ approximates $\mathbf{X}_t^\top \mathbf{X}_t$ well, we could use $\mathbf{S}_t^\top \mathbf{S}_t + \lambda \mathbf{I}$ in place of \mathbf{V}_t . In the following we use the notation $\tilde{\mathbf{V}}_t = \mathbf{S}_t^\top \mathbf{S}_t + \lambda \mathbf{I}$ to denote the sketched regularized correlation matrix. The RLS estimate based upon it is denoted by

$$\tilde{\mathbf{w}}_t = \tilde{\mathbf{V}}_t^{-1} \sum_{s=1}^t \mathbf{x}_s Y_s. \quad (7)$$

A trivial replacement of \mathbf{V} with $\tilde{\mathbf{V}}$ does not yield an efficient algorithm. On the other hand, using the Woodbury identity we may write

$$\tilde{\mathbf{V}}_t^{-1} = \frac{1}{\lambda} \left(\mathbf{I}_{d \times d} - \mathbf{S}_t^\top \mathbf{H}_t \mathbf{S}_t \right)$$

where $\mathbf{H}_t = \left(\mathbf{S}_t \mathbf{S}_t^\top + \lambda \mathbf{I}_{m \times m} \right)^{-1}$. Here matrix-vector multiplications involving \mathbf{S}_t require time $\mathcal{O}(md)$, while matrix-matrix multiplications involving \mathbf{H}_t require time $\mathcal{O}(m^2)$. So, as long as \mathbf{S}_t and \mathbf{H}_t can be efficiently maintained, we obtain an algorithm for linear stochastic bandits where $\tilde{\mathbf{V}}_t^{-1}$ can be updated in time $\mathcal{O}(md + m^2)$. Next, we focus on a concrete sketching algorithm that ensures efficient updates of \mathbf{S}_t and \mathbf{H}_t .

Frequent Directions. Frequent Directions (FD) (Ghashami et al., 2016) is a deterministic sketching algorithm that maintains a matrix \mathbf{S}_t whose last row is invariably $\mathbf{0}$. On each round, we insert \mathbf{x}_t^\top into the last row of \mathbf{S}_{t-1} , perform an eigendecomposition $\mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} + \mathbf{x}_t \mathbf{x}_t^\top = \mathbf{U}_t \boldsymbol{\Sigma}_t \mathbf{U}_t^\top$, and then set $\mathbf{S}_t = \left(\boldsymbol{\Sigma}_t - \rho_t \mathbf{I}_{m \times m} \right)^{\frac{1}{2}} \mathbf{U}_t$, where ρ_t is the smallest eigenvalue of $\mathbf{S}_t^\top \mathbf{S}_t$. Observe that the rows of \mathbf{S}_t form an orthogonal basis, and therefore \mathbf{H}_t is a diagonal matrix which can be updated and stored efficiently. Now, the only step in question is an eigendecomposition, which can also be done in time $\mathcal{O}(md)$ —see (Ghashami et al., 2016, Section 3.2). Hence, the total update time per round is $\mathcal{O}(md)$. The updates of matrices \mathbf{S}_t and \mathbf{H}_t are summarized in Algorithm 4.

Algorithm 4 (FD Sketching)

Input: $\mathbf{S}_{t-1} \in \mathbb{R}^{m \times d}$, $\mathbf{x}_t \in \mathbb{R}^d$, $\lambda > 0$

- 1: Compute eigendecomposition $\mathbf{U}^\top \text{diag}\{\rho_1, \dots, \rho_m\} \mathbf{U} = \mathbf{S}_{t-1}^\top \mathbf{S}_{t-1} + \mathbf{x}_t \mathbf{x}_t^\top$
- 2: $\mathbf{S}_t \leftarrow \text{diag}\{\sqrt{\rho_1 - \rho_m}, \dots, \sqrt{\rho_{m-1} - \rho_m}, 0\} \mathbf{U}$
- 3: $\mathbf{H}_t \leftarrow \text{diag}\left\{\frac{1}{\rho_1 - \rho_m + \lambda}, \dots, \frac{1}{\lambda}\right\}$

Output: $\mathbf{S}_t, \mathbf{H}_t$

It is not hard to see that FD sketching sequentially identifies the top- m eigenvectors of the matrix $\mathbf{X}_T^\top \mathbf{X}_T$. Thus, whenever we use a sketched estimate, we lose a part of the spectrum tail. This loss is captured by the following notion of *spectral error*,

$$\varepsilon_m = \min_{k=0, \dots, m-1} \frac{\lambda_{d-k} + \lambda_{d-k+1} + \dots + \lambda_d}{\lambda(m-k)} \quad (8)$$

where $\lambda_1 \geq \dots \geq \lambda_d$ are the eigenvalues of the correlation matrix $\mathbf{X}_T^\top \mathbf{X}_T$. Note that $\varepsilon_m \leq (\lambda_m + \dots + \lambda_d)/\lambda$. For matrices with low rank or light-tailed spectra we expect this spectral error to be small. In the following, we use \tilde{m} to denote the quantity $m + d \ln(1 + \varepsilon_m)$ which occurs often in our bounds involving sketching. Note that $\tilde{m} \geq m$ and $\tilde{m} \rightarrow m$ as the spectral error vanishes.

Since the matrix \mathbf{V}_t is used to compute both the RLS estimate $\hat{\mathbf{w}}_t$ and the norm $\|\cdot\|_{\mathbf{V}_t}$, the sketching of \mathbf{V}_t clearly affects the confidence ellipsoid. The next theorem quantifies how much the confidence ellipsoid must be blown up in order to compensate for the sketching error. Let ρ_t be the smallest eigenvalue of the FD-sketched correlation matrix $\mathbf{S}_t^\top \mathbf{S}_t$ and let $\bar{\rho}_t = \rho_1 + \dots + \rho_t$. The following proposition due to Ghashami et al. (2016) (see the proof of Thm. 3.1, bound on Δ) relates $\bar{\rho}_t$ to ε_m defined in (8).

Proposition 1. *For any $t = 0, \dots, T$, any $\lambda > 0$, and any sketch size $m = 1, \dots, d$, it holds that $\bar{\rho}_t/\lambda \leq \varepsilon_m$.*

A key lemma in the analysis of regret is the following sketched version of (Abbasi-Yadkori et al., 2011, Lemma 11), which bounds the sum of the ridge leverage scores. Although sketching introduces the spectral error ε_m , it also improves the dependence on the dimension from d to m whenever ε_m is sufficiently small.

Lemma 1 (Sketched leverage scores).

$$\begin{aligned} & \sum_{t=1}^T \min \left\{ 1, \|\mathbf{x}_t\|_{\tilde{\mathbf{V}}_{t-1}}^2 \right\} \\ & \leq 2(1 + \varepsilon_m) \left(\tilde{m} + m \ln \left(1 + \frac{TL^2}{m\lambda} \right) \right). \end{aligned} \quad (9)$$

We can now state the main result of this section.

Theorem 2 (Sketched confidence ellipsoid). *Let $\tilde{\mathbf{w}}_t$ be the RLS estimate constructed by an arbitrary policy for linear contextual bandits after t rounds of play. For any $\delta \in (0, 1)$, the optimal parameter \mathbf{w}^* belongs to the set $\tilde{\mathcal{C}}_t \equiv \left\{ \mathbf{w} \in \mathbb{R}^d : \|\mathbf{w} - \tilde{\mathbf{w}}_t\|_{\tilde{\mathbf{V}}_t} \leq \tilde{\beta}_t(\delta) \right\}$ with probability at least $1 - \delta$, where*

$$\begin{aligned} \tilde{\beta}_t(\delta) = & R \sqrt{m \ln \left(1 + \frac{tL^2}{m\lambda} \right) + 2 \ln \frac{1}{\delta} + d \ln \left(1 + \frac{\bar{\rho}_t}{\lambda} \right)} \\ & \cdot \sqrt{1 + \frac{\bar{\rho}_t}{\lambda}} + S\sqrt{\lambda} \left(1 + \frac{\bar{\rho}_t}{\lambda} \right) \end{aligned} \quad (10)$$

$$\stackrel{\circ}{=} R\sqrt{\tilde{m}(1 + \varepsilon_m)} + S\sqrt{\lambda}(1 + \varepsilon_m). \quad (11)$$

Note that (11) is larger than its non-sketched counterpart (6) due to the factors $1 + \varepsilon_m$. However, when the spectral error ε_m vanishes, $\tilde{\beta}_t(\delta)$ becomes of order $R\sqrt{m} + S\sqrt{\lambda}$, which improves upon (6) since we replace the dependence on the ambient space dimension d with the dependence on the sketch size m . In the following, we use the abbreviation $M_\lambda = \max\{1, 1/\sqrt{\lambda}\}$.

4 Sketched OFUL

Equipped with the sketched confidence ellipsoid and the sketched RLS estimate, we can now introduce SOFUL (Algorithm 5), the sketched version of OFUL. SOFUL

Algorithm 5 (SOFUL)

Input: $\delta, \lambda > 0, m \in \{1, \dots, d-1\}$

- 1: $\tilde{\mathbf{w}}_0 = \mathbf{0}, \tilde{\mathbf{V}}_0^{-1} = \frac{1}{\lambda} \mathbf{I}_{d \times d}, \mathbf{S}_0 = \mathbf{0}_{m \times d}$
- 2: **for** $t = 1, 2, \dots$ **do**
- 3: Get decision set D_t
- 4: Play $\mathbf{x}_t \leftarrow \arg \max_{\mathbf{x} \in D_t} \left\{ \tilde{\mathbf{w}}_{t-1}^\top \mathbf{x} + \tilde{\beta}_{t-1}(\delta) \|\mathbf{x}\|_{\tilde{\mathbf{V}}_{t-1}^{-1}} \right\}$
- 5: Observe reward Y_t
- 6: Compute $\mathbf{S}_t, \mathbf{H}_t$ using Alg. 4 given $\mathbf{S}_{t-1}, \mathbf{x}_t$
- 7: $\tilde{\mathbf{V}}_t^{-1} \leftarrow \frac{1}{\lambda} \left(\mathbf{I}_{d \times d} - \mathbf{S}_t^\top \mathbf{H}_t \mathbf{S}_t \right)$
- 8: Compute $\tilde{\mathbf{w}}_t$ using (7)
- 9: **end for**

enjoys the following regret bound, characterized in terms of the spectral error.

Theorem 3. *The regret of SOFUL with FD-sketching of size m w.h.p. satisfies*

$$R_T \stackrel{\tilde{\mathcal{O}}}{=} M_\lambda (1 + \varepsilon_m)^{\frac{3}{2}} \tilde{m} \left(R + S\sqrt{\lambda} \right) \sqrt{T}.$$

Similarly to Abbasi-Yadkori et al. (2011), we also prove a distribution dependent regret bound for SOFUL. This bound is polylogarithmic in time and depends on the smallest difference Δ between the rewards of the best and the second best action in the decision sets,

$$\Delta = \min_{t=1, \dots, T} \max_{\mathbf{x} \in D_t \setminus \{\mathbf{x}_t^*\}} (\mathbf{x}_t^* - \mathbf{x})^\top \mathbf{w}^*.$$

Theorem 4. *The regret of SOFUL with FD-sketching of size m w.h.p. satisfies*

$$R_T \stackrel{\tilde{\mathcal{O}}}{=} M_\lambda (1 + \varepsilon_m)^3 \tilde{m}^2 (R^2 + S^2 \lambda) \frac{(\ln T)^2}{\Delta}.$$

Proofs of the regret bounds appear in the supplementary material (Section A.3).

5 Sketched linear TS

In this section we introduce a variant of linear TS (Algorithm 3) based on FD-sketching. Similarly to

SOFUL, sketched linear TS (see Algorithm 6) uses the FD-sketched approximation $\tilde{\mathbf{V}}_{t-1}$ of the correlation matrix \mathbf{V}_{t-1} in order to select the action \mathbf{x}_t . Note

Algorithm 6 (Sketched linear TS)

Input: $\delta, \lambda > 0, m \in \{1, \dots, d-1\}, \mathcal{D}^{\text{TS}}$ (TS-sampling distribution)

- 1: $\tilde{\mathbf{w}}_0 = \mathbf{0}, \tilde{\mathbf{V}}_0^{-1} = \frac{1}{\lambda} \mathbf{I}_{d \times d}, \mathbf{S}_0 = \mathbf{0}_{m \times d}, \delta' = \delta/(4T)$
- 2: **for** $t = 1, 2, \dots$ **do**
- 3: Get decision set D_t
- 4: Sample $\mathbf{Z}_t \sim \mathcal{D}^{\text{TS}}$
- 5: Play $\mathbf{x}_t \leftarrow \arg \max_{\mathbf{x} \in D_t} \mathbf{x}^\top \left(\tilde{\mathbf{w}}_{t-1} + \tilde{\beta}_t(\delta') \tilde{\mathbf{V}}_{t-1}^{-\frac{1}{2}} \mathbf{Z}_t \right)$
- 6: Observe reward Y_t
- 7: Compute $\mathbf{S}_t, \mathbf{H}_t$ using Alg. 4 given $\mathbf{S}_{t-1}, \mathbf{x}_t$
- 8: $\tilde{\mathbf{V}}_t^{-1} \leftarrow \frac{1}{\lambda} \left(\mathbf{I}_{d \times d} - \mathbf{S}_t^\top \mathbf{H}_t \mathbf{S}_t \right)$
- 9: Compute $\tilde{\mathbf{w}}_t$ using (7)
- 10: **end for**

that, in this case, we need both $\tilde{\mathbf{V}}_{t-1}^{-1}$ and $\tilde{\mathbf{V}}_{t-1}^{-\frac{1}{2}}$ to compute \mathbf{x}_t . Using the generalized Woodbury identity (Corollary 1 in Appendix A.2 for proofs), we can write

$$\tilde{\mathbf{V}}_t^{-\frac{1}{2}} = \mathbf{S}_t'^\top \left(\mathbf{S}_t' \mathbf{S}_t'^\top \right)^{-1} \left(\frac{\lambda}{2} \mathbf{I} + \mathbf{S}_t' \mathbf{S}_t'^\top \right)^{-\frac{1}{2}} \mathbf{S}_t'$$

where $\mathbf{S}_t' = \left(\boldsymbol{\Sigma}_t + \left(\frac{\lambda}{2} - \rho_t \right) \mathbf{I}_{m \times m} \right)^{\frac{1}{2}} \mathbf{U}_t$. Note that $\tilde{\mathbf{V}}_t^{-\frac{1}{2}}$ can still be computed in time $\mathcal{O}(md + m^2)$ because $\mathbf{S}_t' \mathbf{S}_t'^\top$ is a diagonal matrix.

The confidence ellipsoid stated in Theorem 2 applies to any contextual bandit policy, and so also to the $\tilde{\mathbf{w}}_t$ constructed by sketched linear TS. However, as shown by Abeille and Lazaric (2017), the analysis needs a confidence ellipsoid larger by a factor equal to the bound on $\|\mathbf{Z}\|$ appearing in the concentration property of the TS-sampling distribution. More precisely, the *TS-confidence ellipsoid* is defined by $\tilde{\mathcal{C}}_t^{\text{TS}} \equiv \left\{ \mathbf{w} \in \mathbb{R}^d : \|\mathbf{w} - \tilde{\mathbf{w}}_t\|_{\tilde{\mathbf{V}}_t} \leq \tilde{\gamma}_t(\delta/(4T)) \right\}$ where

$$\tilde{\gamma}_t(\delta) = \tilde{\beta}_t(\delta) \sqrt{cd \ln \left(\frac{c'd}{\delta} \right)}. \quad (12)$$

The quantity $\tilde{\beta}_t(\delta)$ is defined in (10) and c, c' are the concentration constants of the TS-sampling distribution (Definition 1). We are now ready to prove a bound on the regret of linear TS with FD-sketching.

Theorem 5. *The regret of FD-sketched linear TS, run with sketch size m w.h.p. satisfies*

$$R_T \stackrel{\tilde{\mathcal{O}}}{=} M_\lambda (1 + \varepsilon_m)^{\frac{3}{2}} \tilde{m} \left(R + S\sqrt{\lambda} \right) \sqrt{dT}.$$

The proof of Theorem 5 closely follows the analysis of Abeille and Lazaric (2017) with some key modifications due to the sketching operations. For completeness, we include the proof in the supplementary material.

6 Some proof sketches

Our regret analyses follow Abbasi-Yadkori et al. (2011); Abeille and Lazaric (2017) and related works. However, due to the sketching of the correlation matrix, some key components of the proofs now depend on the spectral error (8). Because of that, we need tools specific to the analysis of linear bandits with FD-sketching. These tools are used to bound the instantaneous regret $(\mathbf{x}^* - \mathbf{x}_t)^\top \mathbf{w}^*$ in terms of the norm $\|\mathbf{w}^* - \tilde{\mathbf{w}}_t\|_{\tilde{\mathbf{V}}_{t-1}}$ and the ridge leverage scores. Armed with these results, we then prove our regret bounds in Section 6.1. In this section we only present the core regret analysis of SOFUL, while we defer its complete analysis and analysis of the sketched linear TS to the supplementary material.

6.1 Proof of regret bounds for SOFUL

We start by introducing a basic relationship between the correlation matrix of actions $\mathbf{X}_s^\top \mathbf{X}_s$ and its FD-sketched estimate $\mathbf{S}_t^\top \mathbf{S}_t$ with sketch size $m \leq d$. Recall that ρ_t is the smallest eigenvalue of $\mathbf{S}_t^\top \mathbf{S}_t$ for $t = 1, \dots, T$ and $\bar{\rho}_t = \rho_1 + \dots + \rho_t$.

Proposition 2. *Let \mathbf{S}_s be the matrix computed by FD-sketching at time step $s = 1, \dots, t$ (where $\mathbf{S}_0 = \mathbf{0}$). Then $\mathbf{X}_s^\top \mathbf{X}_s = \mathbf{S}_s^\top \mathbf{S}_s + \bar{\rho}_s \mathbf{I}$.*

Proof. By construction, $\mathbf{S}_{s-1}^\top \mathbf{S}_{s-1} + \mathbf{x}_s \mathbf{x}_s^\top = \mathbf{U}_s \boldsymbol{\Sigma}_s \mathbf{U}_s^\top$ where $\mathbf{S}_s = (\boldsymbol{\Sigma}_s - \rho_s \mathbf{I}_{m \times m})^{\frac{1}{2}} \mathbf{U}_s$. Thus,

$$\mathbf{S}_s^\top \mathbf{S}_s = \mathbf{U}_s \boldsymbol{\Sigma}_s \mathbf{U}_s^\top - \rho_s \mathbf{I} = \mathbf{S}_{s-1}^\top \mathbf{S}_{s-1} + \mathbf{x}_s \mathbf{x}_s^\top - \rho_s \mathbf{I}$$

Summing both sides of the above over $s = 1, \dots, t$, $\mathbf{S}_t^\top \mathbf{S}_t = \sum_{s=1}^t \mathbf{x}_s \mathbf{x}_s^\top - \sum_{s=1}^t \rho_s \mathbf{I}$. \square

We will also use the following lemma of (Abbasi-Yadkori et al., 2011, Lemma 11).

Lemma 2. *For $\lambda \geq \max\{1, L^2\}$, we have that*

$$\sum_{t=1}^T \|\mathbf{x}_t\|_{\mathbf{V}_{t-1}}^2 \leq 2d \ln \left(1 + \frac{TL^2}{\lambda d} \right). \quad (13)$$

The following lemma gives a sketch-specific version of the determinant-trace inequality, shown in the supplementary material (Lemma 8).

Lemma 3.

$$\ln \left(\frac{\det(\mathbf{V}_t)}{\det(\lambda \mathbf{I})} \right) \leq d \ln \left(1 + \frac{\bar{\rho}}{\lambda} \right) + m \ln \left(1 + \frac{tL^2}{m\lambda} \right).$$

Next we prove Lemma 1, which is similar to Lemma 2. However, now the statement depends on the sketched matrix $\tilde{\mathbf{V}}_{t-1}$ instead of \mathbf{V}_{t-1} .

Proof of Lemma 1. Throughout the proof, unless stated explicitly, we drop the subscripts containing t . Therefore, $\mathbf{V} = \mathbf{V}_{t-1}$, $\tilde{\mathbf{V}} = \tilde{\mathbf{V}}_{t-1}$, $\mathbf{x} = \mathbf{x}_t$, and $\bar{\rho} = \bar{\rho}_{t-1}$. Now suppose that $(\tilde{\lambda}_i + \lambda, \tilde{\mathbf{u}}_i)$ is an i -th eigenpair of $\tilde{\mathbf{V}}$. Then, Proposition 2 implies that a corresponding eigenpair of \mathbf{V} is $(\tilde{\lambda}_i + \lambda + \bar{\rho}, \tilde{\mathbf{u}}_i)$. Using this fact we have that

$$\begin{aligned} \|\mathbf{x}\|_{\mathbf{V}^{-1}}^2 &= \mathbf{x}^\top \tilde{\mathbf{V}} \tilde{\mathbf{V}}^{-1} \mathbf{V}^{-1} \mathbf{x} \\ &= \mathbf{x}^\top \left(\sum_{i=1}^d \tilde{\mathbf{u}}_i \tilde{\mathbf{u}}_i^\top \frac{1}{\tilde{\lambda}_i + \lambda} \frac{\tilde{\lambda}_i + \lambda}{\tilde{\lambda}_i + \lambda + \bar{\rho}} \right) \mathbf{x} \\ &\geq \frac{\lambda}{\lambda + \bar{\rho}} \mathbf{x}^\top \left(\sum_{i=1}^d \tilde{\mathbf{u}}_i \tilde{\mathbf{u}}_i^\top \frac{1}{\tilde{\lambda}_i + \lambda} \right) \mathbf{x} = \frac{\lambda}{\lambda + \bar{\rho}} \|\mathbf{x}\|_{\tilde{\mathbf{V}}^{-1}}^2. \end{aligned}$$

Furthermore, this implies that

$$\begin{aligned} \min \left\{ 1, \frac{\lambda}{\lambda + \bar{\rho}} \|\mathbf{x}\|_{\tilde{\mathbf{V}}^{-1}}^2 \right\} &\leq \min \{ 1, \|\mathbf{x}\|_{\mathbf{V}^{-1}}^2 \} \\ \Rightarrow \min \left\{ 1, \|\mathbf{x}\|_{\tilde{\mathbf{V}}^{-1}}^2 \right\} &\leq \left(1 + \frac{\bar{\rho}}{\lambda} \right) \min \{ 1, \|\mathbf{x}\|_{\mathbf{V}^{-1}}^2 \}. \end{aligned}$$

Finally, combining the above with Lemma 2, equation and using the fact that $\bar{\rho}_{t-1} \leq \bar{\rho}_T$, we obtain

$$\begin{aligned} \sum_{t=1}^T \min \left\{ 1, \|\mathbf{x}_t\|_{\tilde{\mathbf{V}}_{t-1}}^2 \right\} &\leq 2 \left(1 + \frac{\bar{\rho}_T}{\lambda} \right) \ln \left(\frac{\det(\mathbf{V}_T)}{\det(\lambda \mathbf{I})} \right) \\ &\leq 2 \left(1 + \frac{\bar{\rho}_T}{\lambda} \right) \left(d \ln \left(1 + \frac{\bar{\rho}_T}{\lambda} \right) + m \ln \left(1 + \frac{TL^2}{m\lambda} \right) \right) \\ &\leq 2(1 + \varepsilon_m) \left(d \ln(1 + \varepsilon_m) + m \ln \left(1 + \frac{TL^2}{m\lambda} \right) \right) \end{aligned}$$

where the penultimate inequality follows from Lemma 3 and the last step follows from Proposition 1. \square

Now we give a bound on the instantaneous regret.

Lemma 4. *For any $\delta > 0$, the instantaneous regret of SOFUL satisfies $(\mathbf{x}_t^* - \mathbf{x}_t)^\top \mathbf{w}^* \leq 2\tilde{\beta}_{t-1}(\delta) \|\mathbf{x}_t\|_{\tilde{\mathbf{V}}_{t-1}^{-1}}$ for $t = 1, \dots, T$.*

Proof. Let $\tilde{\mathbf{w}}_{t-1}^{\text{SO}}$ be the FD-sketched RLS estimate of OFUL (Algorithm 5). Recall that the optimal action at time t is $\mathbf{x}_t^* = \arg \max_{\mathbf{x} \in D_t} \mathbf{x}^\top \mathbf{w}^*$, whereas

$$(\mathbf{x}_t, \tilde{\mathbf{w}}_{t-1}^{\text{SO}}) = \arg \max_{(\mathbf{x}, \mathbf{w}) \in D_t \times \tilde{C}_{t-1}} \mathbf{x}^\top \mathbf{w}$$

We use these facts to bound the instantaneous regret,

$$\begin{aligned} (\mathbf{x}_t^* - \mathbf{x}_t)^\top \mathbf{w}^* &\leq \mathbf{x}_t^\top (\tilde{\mathbf{w}}_{t-1}^{\text{SO}} - \mathbf{w}^*) \\ &= \mathbf{x}_t^\top (\tilde{\mathbf{w}}_{t-1}^{\text{SO}} - \tilde{\mathbf{w}}_{t-1}) + \mathbf{x}_t^\top (\tilde{\mathbf{w}}_{t-1} - \mathbf{w}^*) \end{aligned}$$

$$\begin{aligned} &\leq \|\mathbf{x}_t\|_{\tilde{\mathbf{V}}_{t-1}^{-1}} \left(\|\tilde{\mathbf{w}}_{t-1}^{\text{SO}} - \tilde{\mathbf{w}}_{t-1}\|_{\tilde{\mathbf{V}}_{t-1}} + \|\tilde{\mathbf{w}}_{t-1} - \mathbf{w}^*\|_{\tilde{\mathbf{V}}_{t-1}} \right) \\ &\leq 2\tilde{\beta}_{t-1}(\delta) \|\mathbf{x}_t\|_{\tilde{\mathbf{V}}_{t-1}^{-1}} \quad (\text{by Theorem 2}) \end{aligned}$$

where we get penultimate inequality by Cauchy-Schwartz inequality. \square

Proof of Theorem 3. Using Lemma 4 gives

$$\begin{aligned} R_T &= \sum_{t=1}^T (\mathbf{x}_t^* - \mathbf{x}_t)^\top \mathbf{w}^* \\ &\leq 2 \sum_{t=1}^T \min \left\{ LS, \tilde{\beta}_{t-1}(\delta) \|\mathbf{x}_t\|_{\tilde{\mathbf{V}}_{t-1}^{-1}} \right\} \quad (14) \end{aligned}$$

$$\begin{aligned} &\leq 2 \sum_{t=1}^T \tilde{\beta}_{t-1}(\delta) \min \left\{ \frac{L}{\sqrt{\lambda}}, \|\mathbf{x}_t\|_{\tilde{\mathbf{V}}_{t-1}^{-1}} \right\} \quad (15) \\ &\leq 2\tilde{A} \sum_{t=1}^T \min \left\{ \frac{L}{\sqrt{\lambda}}, \|\mathbf{x}_t\|_{\tilde{\mathbf{V}}_{t-1}^{-1}} \right\} \\ &\leq 2\tilde{A} \max \left\{ 1, \frac{L}{\sqrt{\lambda}} \right\} \sum_{t=1}^T \min \left\{ 1, \|\mathbf{x}_t\|_{\tilde{\mathbf{V}}_{t-1}^{-1}} \right\} \\ &\leq 2\tilde{A} \max \left\{ 1, \frac{L}{\sqrt{\lambda}} \right\} \sqrt{T \sum_{t=1}^T \min \left\{ 1, \|\mathbf{x}_t\|_{\tilde{\mathbf{V}}_{t-1}^{-1}}^2 \right\}} \end{aligned}$$

where we used $\tilde{A} = \max_{t=0, \dots, T-1} \tilde{\beta}_t(\delta)$. Also, we get (14) since $\max_{t=1, \dots, T} \max_{\mathbf{x} \in D_t} |\mathbf{x}^\top \mathbf{w}^*| \leq LS$ by Cauchy-Schwartz, and (15) since $\min_{t=0, \dots, T-1} \min_{\delta \in [0, 1]} \tilde{\beta}_t(\delta) \geq S\sqrt{\lambda}$. The last inequality is obtained via the Cauchy-Schwartz inequality. Now we finish by bounding $\tilde{\beta}_t(\delta)$ using (11) and we bound the summation term using Lemma 1,

$$\begin{aligned} R_T &\stackrel{\tilde{\mathcal{O}}}{=} M_\lambda \sqrt{T} \left(R\sqrt{\tilde{m}(1 + \varepsilon_m)} + S\sqrt{\lambda}(1 + \varepsilon_m) \right) \\ &\quad \cdot \sqrt{\tilde{m}(1 + \varepsilon_m)} \\ &\stackrel{\tilde{\mathcal{O}}}{=} M_\lambda \sqrt{T} \left(R\tilde{m}(1 + \varepsilon_m) + S\sqrt{\lambda}(1 + \varepsilon_m)^{\frac{3}{2}} \sqrt{\tilde{m}} \right) \\ &\stackrel{\tilde{\mathcal{O}}}{=} M_\lambda (1 + \varepsilon_m)^{\frac{3}{2}} \tilde{m} \left(R + S\sqrt{\lambda} \right) \sqrt{T} \quad \square \end{aligned}$$

7 Experiments

In this section we present a simple empirical evaluation of OFUL and linear TS against their sketched versions.

Setup. The idea of our experimental setup is similar to the one described by Cesa-Bianchi et al. (2013). We convert a K -class classification problem into a contextual bandit problem as follows: given a dataset of labeled instances $(\mathbf{x}, y) \in \mathbb{R}^d \times \{1, \dots, K\}$, we partition it into K subsets according to the class labels. Then we

create K sequences by drawing a random permutation of each subset. At each step t the decision set D_t is obtained by picking the t -th instance from each one of these K sequences. Finally, rewards are determined by choosing a class $y \in \{1, \dots, K\}$ and then consistently assigning reward 1 to all instances labeled with y and reward 0 to all remaining instances. The reported mean and standard deviation of the cumulative reward are averaged over four random permutations of the K sequences. This procedure is used for all baselines.

Datasets. We perform experiments on six publicly available datasets for multiclass classification from the `openml` repository (Vanschoren et al., 2013), see the table below here for details.

Dataset	Examples	Features	Classes
Bank	45k	17	2
SatImage	6k	37	6
Spam	4k	58	2
Pendigits	11k	17	10
MFeat	2k	48	10
CMC	1.4k	10	3

Baselines. The hyperparameters β (confidence ellipsoid radius) and λ (RLS regularization parameter) are selected on a validation set of size 100 via grid search on $(\beta, \lambda) \in \{1, 10^2, 10^3, 10^4\} \times \{10^{-2}, 10^{-1}, 1\}$ for OFUL, and $\{1, 10^2, 10^3\} \times \{10^{-2}, 10^{-1}, 1, 10^2\}$ for linear TS.

Results The experiments we present in this section measure cumulative reward throughout time. In the first experiment we compare OFUL and linear TS to their sketched versions while varying the sketch size (in the plots the sketch size is expressed as a percentage of the context space dimension). Results for the three datasets are presented in Figure 1, while results for the remaining datasets can be found in the supplementary material. We observe that on two datasets out of three, sketched algorithms indeed do not suffer a substantial drop in performance when compared to the non-sketched ones, even when the sketch size amounts to 60% of the context space dimension. This demonstrates that sketching successfully captures relevant subspace information relatively to the goal of maximizing reward.

Because the FD-sketching procedure considered in this paper is essentially performing online PCA, it is natural to ask how our sketched algorithms would compare to their non-sketched version run on the best m -dimensional subspace (computed by running PCA on the entire dataset). In Figure 2, we compare OFUL and SOFUL (results for linear TS are included in the supplementary material). In particular, we keep 60%, 40%, and 20% of the top principal components, and notice that, like in Figure 1, there are cases with little or no loss in performance.

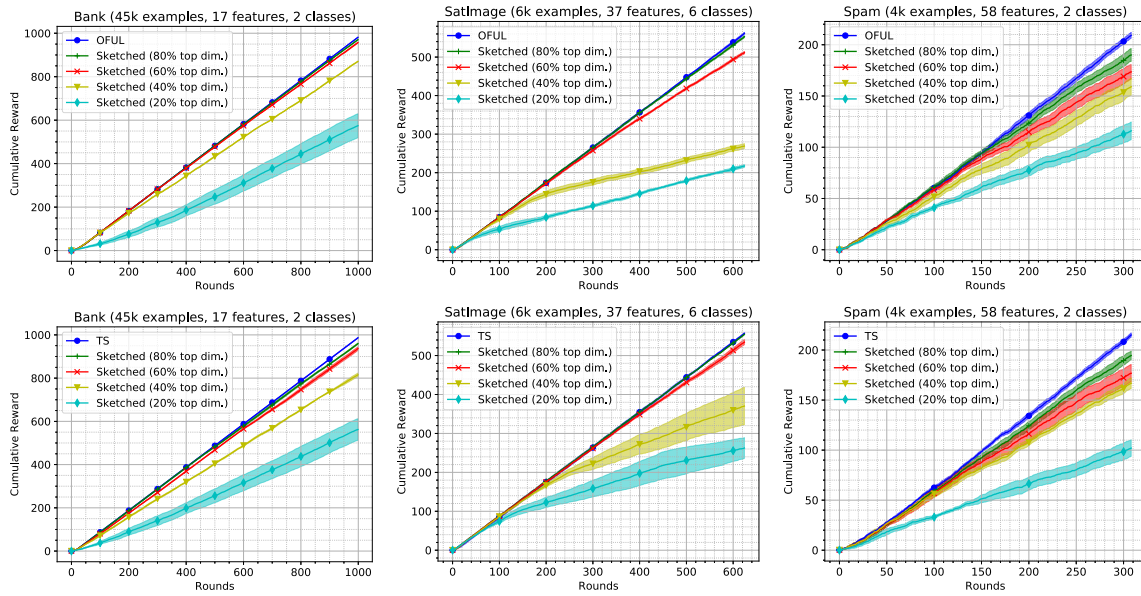


Figure 1: Comparison of SOFUL (first row) and sketched linear TS (second row) to their non-sketched variants on three real-world datasets and for different sketch sizes. Note that, in some cases, a sketch size equal to 80% and even 60% of the context space dimension does not significantly affect the performance.

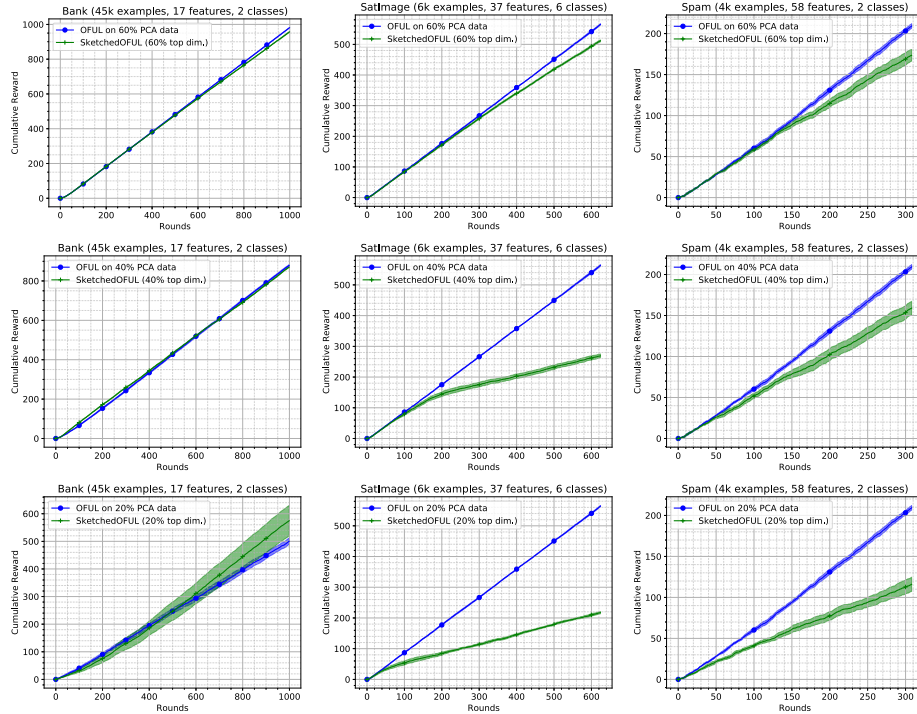


Figure 2: Comparison of OFUL run on the best m -dimensional subspace against SOFUL run with sketch size m . Rows show m as a fraction of the context space dimension: 60%, 40%, 20%, while columns correspond to different datasets. Note that, in some cases (with sketch size m of size at least 60%), SOFUL performs as well as if the best m -dimensional subspace had been known in hindsight.

References

- Y. Abbasi-Yadkori, D. Pál, and C. Szepesvári. Improved algorithms for linear stochastic bandits. In *Conference on Neural Information Processing Systems (NIPS)*, pages 2312–2320, 2011.
- M. Abeille and A. Lazaric. Linear Thompson sampling revisited. *Electronic Journal of Statistics*, 11(2):5165–5197, 2017.
- S. Agrawal and N. Goyal. Thompson sampling for contextual bandits with linear payoffs. In *International Conference on Machine Learning (ICML)*, pages 127–135, 2013.
- P. Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research*, 3(Nov):397–422, 2002.
- D. Calandriello, A. Lazaric, and M. Valko. Efficient second-order online kernel learning with adaptive embedding. In *Conference on Neural Information Processing Systems (NIPS)*, pages 6140–6150, 2017.
- N. Cesa-Bianchi, C. Gentile, and G. Zappella. A gang of bandits. In *Conference on Neural Information Processing Systems (NIPS)*, pages 737–745, 2013.
- V. Dani, T. P. Hayes, and S. M. Kakade. Stochastic linear optimization under bandit feedback. In *Conference on Computational Learning Theory (COLT)*, pages 355–366, 2008.
- M. Ghashami, E. Liberty, J. M. Phillips, and D. P. Woodruff. Frequent directions: Simple and deterministic matrix sketching. *SIAM Journal on Computing*, 45(5):1762–1792, 2016.
- A. Gonen, F. Orabona, and S. Shalev-Shwartz. Solving ridge regression using sketched preconditioned SVRG. In *International Conference on Machine Learning (ICML)*, pages 1397–1405, 2016.
- N. J. Higham. *Functions of matrices: theory and computation*, volume 104. Siam, 2008.
- K.-S. Jun, A. Bhargava, R. Nowak, and R. Willett. Scalable generalized linear bandits: Online computation and hashing. In *Conference on Neural Information Processing Systems (NIPS)*, pages 99–109, 2017.
- T. Lattimore and C. Szepesvári. *Bandit Algorithms*. Cambridge University Press, 2018.
- H. Luo, A. Agarwal, N. Cesa-Bianchi, and J. Langford. Efficient second order online learning by sketching. In *Conference on Neural Information Processing Systems (NIPS)*, pages 902–910, 2016.
- J. Vanschoren, J. N. van Rijn, B. Bischl, and L. Torgo. OpenML: Networked Science in Machine Learning. *SIGKDD Explorations*, 15(2):49–60, 2013.
- D. Woodruff. Sketching as a tool for numerical linear algebra. *Foundations and Trends in Theoretical Computer Science*, 10(1–2):1–157, 2014.
- X. Yu, M. R. Lyu, and I. King. Cbrap: Contextual bandits with random projection. In *Conference on Artificial Intelligence (AAAI)*, pages 2859–2866, 2017.