# Stochastic Online Linear Regression: the Forward Algorithm to Replace Ridge

**Reda Ouhamma**
Univ. Lille, CNRS, Inria, Centrale Lille, UMR 9189 - CRIStAL, F-59000
`reda.ouhamma@univ-lille.fr`

**Odalric. Maillard**
Univ. Lille, CNRS, Inria, Centrale Lille, UMR 9189 - CRIStAL, F-59000

**Vianney. Perchet**
Criteo, ENSAE, ENS PARIS-SACLAY

## Abstract

We consider the problem of online linear regression in the stochastic setting. We derive high probability regret bounds for online *ridge* regression and the *forward* algorithm. This enables us to compare online regression algorithms more accurately and eliminate assumptions of bounded observations and predictions. Our study advocates for the use of the forward algorithm in lieu of ridge due to its enhanced bounds and robustness to the regularization parameter. Moreover, we explain how to integrate it in algorithms involving linear function approximation to remove a boundedness assumption without deteriorating theoretical bounds. We showcase this modification in linear bandit settings where it yields improved regret bounds. Last, we provide numerical experiments to illustrate our results and endorse our intuitions.

## 1 Introduction and preliminaries

The *forward regression* algorithm, popularized in [24, 3], shows competitive performance bounds in the challenging setup of online regression with *adversarial bounded* observations. We revisit the analysis of this strategy in the practically relevant alternative situation of *stochastic* linear regression with sub-Gaussian noise, hence possibly *unbounded* observations. When compared to the classical ridge regression strategy - its natural competitor - the existing analysis in the adversarial bounded case suggests the forward algorithm has higher performances. It is then natural to ask whether this conclusion holds for the stochastic setup. However, we show that in the stochastic setup, the existing adversarial analysis does not seem sufficient to draw conclusions, as it does not capture some important phenomena, such as the concentration of the parameter estimate around the regression parameter. It may further lead the practitioner to use an improper tuning of the regularization parameter. In order to overcome these issues, we revisit the analysis of the forward algorithm in the case of unbounded sub-Gaussian linear regression and provide a high probability regret bound on the performance of the forward and ridge regression strategies. Owing to this refined analysis, we show that the forward algorithm is superior in this scenario as well, but for different reasons than what is suggested by the adversarial analysis. We discuss the implications of this result in a practical application: stochastic linear bandits, both from theoretical and experimental perspectives.

**Setup:** In the classical setting of online regression with the square loss, an environment initially chooses a sequence of feature vectors $\{x_t\}_t \in \mathbb{R}^d$ together with corresponding observations $\{y_t\}_t \in$

$\mathbb{R}$. Then, at each decision step $t$, the learner receives feature vector $x_t$ and must output a prediction $\hat{y}_t \in \mathbb{R}$. Afterwards, the environment reveals the true label $y_t$ and iteration $t+1$ begins. In this article, we focus on the case when the data generating process is a *stochastic* linear model:

$$\exists \theta_* \in \mathbb{R}^d \text{ such that } \forall t \in \mathbb{N}^* : \quad y_t = x_t^\top \theta_* + \epsilon_t,$$

where $\{\epsilon_t\}_t$ is a noise sequence. At iteration $t$, strategy $\mathcal{A}$ computes a parameter $\theta_{t-1}^\mathcal{A}$ to predict $\hat{y}_t^\mathcal{A} = x_t^\top \theta_{t-1}^\mathcal{A}$. In the sequel, we omit the subscript $\mathcal{A}$ when the algorithm is clear from context. The learner's prediction incurs the loss: $\ell_t^\mathcal{A} \overset{\text{def}}{=} \ell(x_t^\top \theta_{t-1}, y_t) = (\hat{y}_t - y_t)^2$, the learner then updates its prediction $\theta_{t-1}$ to $\theta_t$ and so on. The total cumulative loss at horizon $T$ is denoted $L_T^\mathcal{A} = \sum_{t=1}^T \ell_t^\mathcal{A}$. We also let $\ell_t(\theta) = \ell(x_t^\top \theta, y_t)$ (resp. $L_T(\theta) = \sum_{t=1}^T \ell_t(\theta)$) be the instantaneous (resp. cumulative) loss incurred by predicting $\theta$ at time $t$ (resp. $\forall t = 1, \ldots, T$). Online regression algorithms are evaluated using different regret definitions, in the form of a relative cumulative loss to a batch loss; The quantity of interest in this paper is:

$$R_T^\mathcal{A} = L_T^\mathcal{A} - \min_\theta L_T(\theta). \tag{1}$$

From the perspective of online learning theory, online regression algorithms are usually designed for an *adversarial* setting, assuming an arbitrary bounded response variable $|y_t| \leq Y$ at each time step. While the mere existence of algorithms with tight guarantees in this general setting is remarkable, a practitioner may also consider alternative settings, in which analysis for the adversarial setup may be overly conservative. For illustration, we focus on the practical setting of bounded parameter $\|\theta_*\|_2 \leq S$ and i.i.d zero-mean $\sigma$-sub-Gaussian noise sequences:

$$\forall t \geq 1, \gamma \in \mathbb{R} : \quad \mathbb{E}\left[\exp(\gamma\epsilon)\right] \leq \exp(\sigma^2\gamma^2/2).$$

We emphasize that while previous results in literature are valid for the adversarial bounded setting, we will still shed new light on the performance of these strategies in a stochastic unbounded setup, which is neither more general nor more restrictive than the adversarial one, and discuss their implications for the practitioner. Let us recall the two popular online regression algorithms considered.

**Online ridge regression [Algorithm 1]:** This folklore algorithm is defined in the online setting as the greedy version of batch ridge regression:

$$\theta_t^r \in \arg\min_\theta L_t(\theta) + \lambda\|\theta\|_2^2, \tag{2}$$

where $\lambda$ is a parameter and $\lambda\|\theta\|_2^2$ is a regularization used to penalize model complexity.

---
**Algorithm 1:** Online ridge regression
---
Given $\theta_0 \in \mathbb{R}^d$
**for** $t = 1, \ldots, T$ **do**
    observe $x_t \in \mathbb{R}^d$ and predict $\hat{y}_t = x_t^\top \theta_{t-1}^r \in \mathbb{R}$
    observe $y_t$ and incur loss $\ell_t \in \mathbb{R}$
    update parameter: $\theta_t^r \in \arg\min_\theta L_t(\theta) + \lambda\|\theta\|_2^2$
**end**
---

A solution to the quadratic optimization problem of Eq. 2 is given in closed form, by $\theta_t^r = G_t(\lambda)^{-1}b_t$, where $G_t(\lambda) = \lambda I + \sum_{q=1}^t x_q x_q^\top$ and $b_t = \sum_{q=1}^t x_q y_q$. We may further denote $G_t$ instead of $G_t(\lambda)$ when $\lambda$ is clear from context.

**The forward algorithm [Algorithm 2]:** A subtle change to the ridge regression takes advantage of the next feature $x_{t+1}$ to better adapt to the next loss:

$$\theta_t^f \in \arg\min_\theta L_t(\theta) + (x_{t+1}^\top \theta)^2 + \lambda\|\theta\|_2^2. \tag{3}$$

**Algorithm 2:** The forward algorithm

---

Given $\theta_0 \in \mathbb{R}^d$
**for** $t = 1, \ldots, T$ **do**

    observe $\mathrm{x}_t \in \mathbb{R}^d$

    update parameter: $\theta_{t-1}^{\mathsf{f}} \in \arg\min_\theta L_{t-1}(\theta) + (x_t^\top \theta)^2 + \lambda ||\theta||_2^2$

    predict $\hat{y}_t = x_t^\top \theta_{t-1}^{\mathsf{f}} \in \mathbb{R}$

    observe $y_t$ and incur loss $\ell_t \in \mathbb{R}$

**end**

---

Equivalently, the update step can be written: $\theta_t^{\mathsf{f}} = G_{t+1}^{-1} b_t$, where $G_t$ is still defined same as before. Intuitively, the term $(x_{t+1}^\top \theta)^2$ in Eq. 3 is a "predictive loss", a penalty on the parameter $\theta$ in the direction of the new feature vector $x_{t+1}$. This approach can be linked to transductive methods for regression [8, 22]. [22] describe two algorithms for linear prediction in supervised settings, and leverage the knowledge of the next test point to improve the prediction accuracy. However, these algorithms have significant computational complexities and are not adapted to online settings.

**Related work**   Linear regression is perhaps one of the most known algorithms in machine learning, due to is simplicity and explicit solution. In contrast with the *batch* setting (when all observations are provided), *online* linear regression started receiving interest relatively recently. The first theoretical analyses date back to [9, 14, 7, 13]. Under the assumption that the response variable is bounded $|y_t| \leq Y$, it has been shown that the forward algorithm [24, 3] achieves a relative cumulative online error of $dY^2 \log(T)$ compared to the best batch regression strategy. This bound holds *uniformly* over bounded response variables and competitor vectors, and is 4 times better than the corresponding bound derived for online *ridge* regression.

Bartlett et al. [4] studied minimax regret bounds for online regression, and ingeniously removed a dependence on the scale of features in existing bounds by considering the beforehand-known features setting, where all feature points $\{x_t\}_{1 \leq t \leq T}$ are known before the learning starts. Moreover, they derive a "backward algorithm" that is optimal under certain intricate assumptions on observations and features. Later on, [16] were able to prove that under new (tricky) assumptions on observed features and labels the *backward algorithm* is not only optimal but applicable in sequential settings as well. More recently, [11] provided an optimal algorithm in the setting of beforehand known features without imposing stringent conditions as in [4, 16]. They show that the forward algorithm with $\lambda = 0$ yields a first-order *optimal* asymptotic regret bound uniform over bounded observations. However, due to the lack of regularization, their bound (*cf.* Theorem 11 in [11]) may blow up if the design matrix $G_t(0)$ is not full rank. It is hence not uniform over all bounded feature sequences $\{x_t\}_t$.

**Paper outline and contributions:**   In this paper, we continue the line of work initiated on the *forward* algorithm and advocate for its use in the stochastic setting with possibly unbounded response variables, in replacement for the ridge regression (whenever possible). To this end, we consider an online *stochastic* linear regression setup where the noise is assumed to be i.i.d $\sigma$-sub-Gaussian.

In Section 2 we recall the online performance bounds established for ridge regression and the forward algorithm in the *adversarial* case with bounded observations. Next, in subsection 2.3, we discuss some limitations of the adversarial results when comparing regression algorithms in the stochastic setting. For instance, these bounds compare the cumulative loss of a strategy to the value of the batch optimization problem, which may not be indicative of the real performance of the strategy (*cf.* Corollary 2.3.1) and may encourage a sub-optimal tuning of the regularization parameter.

In Section 3, we study the performance of these algorithms using the cumulative regret with respect to the true parameter (*cf.* Eq. 6), which we believe is more practitioner-friendly than comparing to the batch optimization problem. We show in Theorem 3.1 how these two measures of performance are related. We provide in Theorems 3.2 and 3.3 a novel analysis of online regression algorithms without assuming bounded observations. This key result is made possible by considering high probability bounds instead of bounded individual sequences. We show that the regret upper-bound for ridge regression is inversely proportional to the regularization parameter. Consequently, we argue that following these results, forward regression should be used *in lieu* of ridge regression.

In Section 4, we revisit the linear bandit setup previously analyzed assuming *bounded* rewards: we relax this assumption and provide an -optimism in the face of uncertainty- style algorithm with the

forward algorithm instead of ridge, which is especially well-suited for the bandit setup, and provide novel regret analysis in Theorem 4.1. We proceed similarly in Appendix F, revisiting a setup of non-stationary (abruptly changing) linear bandits.

## 2 Adversarial bounds and limitations

In this section, we recall existing results regarding the aforementioned ridge and forward algorithms. We then discuss their limits and benefits when considered from a stochastic perspective.

### 2.1 Adversarial regret bounds (existing results)

One of the first theoretical analyses of online regression dates back to [24] and [3], and is recalled in the theorem below. It is stated in the form of an "online-to-offline conversion" performance bound.

**Theorem 2.1.** *(Theorem 4.6[1] of [3]) The online ridge regression algorithm satisfies:*

$$L_T^r - \min_\theta \left( L_T(\theta) + \lambda \|\theta\|_2^2 \right) \leq 4(Y^r)^2 d \log \left( 1 + \frac{TX^2}{\lambda d} \right),$$

*where* $X = \max_{1 \leq t \leq T} \|x_t\|_2$, *and* $Y^r = \max_{1 \leq t \leq T} \left\{ |y_t|, \left| \boldsymbol{x}_t^\top \boldsymbol{\theta}_{t-1} \right| \right\}$.

The reader should note that this result compares the learner's online cumulative loss to the regularized batch ridge regression loss. As such, it is an online-to-offline conversion regret. This is different from the sequential regret that would compare to the minimum achievable loss. This theorem highlights a dependence on the *range* of predictions of the algorithm, as $Y^r \geq \max_{1 \leq t \leq T} \left| \boldsymbol{x}_t^\top \boldsymbol{\theta}_{t-1} \right|$.

**Remark 1.** *(Small losses) The regret bound for ridge regression can be improved if the learner knows that the loss is small for the best expert, see Orabona et al. [18]. Note however that such techniques require prior knowledge of all the best expert loss $L_T^*$, their optimal bound is* $\sim O(\sqrt{L_T^* \log T})$.

The forward algorithm has an enhanced performance in this setup according to this next result.

**Theorem 2.2.** *(Theorem 5.6 of [3]) The forward algorithm satisfies[2]:*

$$L_T^f - \min_\theta \left( L_T(\theta) + \lambda \|\theta\|_2^2 \right) \leq (Y^f)^2 d \log \left( 1 + \frac{TX^2}{\lambda d} \right),$$

*where* $X = \max_{1 \leq t \leq T} \|x_t\|_2$, *and* $Y^f = \max_{1 \leq t \leq T} |y_t|$.

Notice that in this result $Y$ is different than in Theorem 2.1 and is independent from the algorithm's predictions. Moreover, Theorem 2.2 exhibits a bound that is at least 4 times better than Theorem 2.1. More precisely, Theorem 2.1 suggests that, in order to compare the two bounds, prior knowledge of $Y^f$ is required to further clip the predictions of online ridge regression in $[-Y^f, Y^f]$; and that even with such knowledge the forward algorithm may be 4-times better than ridge regression. We believe that this unfortunately led researchers to turn away from analyzing more deeply what may happen.

### 2.2 Limitation in the adversarial setup: rigid regularization

To evaluate online regression strategies, a tight *lower* bound was derived in [11]. The latter studied uniform minimax lower bounds in the setting of beforehand-known features (that is when $(x_t)_{1 \leq t \leq T}$ known in advance), which is very challenging for a lower bound. They show that, the minimax uniform regret bound is controlled as follows.

**Theorem 2.3.** *(Gaillard et al. [11]) For all $T \geq 8, Y > 0$ we have:*

$$R_{T,[-Y,Y]}^\star \geq dY^2 (\log(T) - (3 + \log(d)) - \log(\log(d))).$$

*where* $R_{T,[-Y,Y]}^\star \overset{def}{=} \inf_{\mathcal{A}} \sup_{x_1,\ldots,_T \in [0,1]^d} \sup_{|y_t| \leq Y} \left\{ \sum_{t=1}^T (y_t - \hat{y}_t^{\mathcal{A}})^2 - \inf_{u \in \mathbb{R}^d} \sum_{t=1}^T (y_t - x_t^\top u)^2 \right\}$

---

[1] Note that there are typos in the statement of this theorem in original paper: compare Lemma 4.2 with Theorems 4.6 and 5.6 therein to see this. Reported theorems are accurate.

[2] See footnote 1.

4

We will use this result to evaluate the optimality of ridge and forward regressions. First, we need to convert Theorems 2.1 and 2.2 to sequential *regret* bounds. Indeed, in their current form, they compare the cumulative loss of the learner to the value of a regularized batch optimization. This next result transforms them, and is a corollary of Theorems 11.7 and 11.8 of Cesa-Bianchi & Lugosi [6].

**Corollary 2.3.1.** *(Of Theorems 11.7 and 11.8 of [6]) For all $T \geq 1$, $(x_t)_{1 \leq t \leq T} \in \mathbb{R}^d$, $(y_t)_{1 \leq t \leq T} \in [-Y, Y]$ such that $\|x_t\|_2 \leq X$,*

$$\text{for } \mathcal{A} \in \{\mathbf{r}, \mathbf{f}\} \qquad R_T^{\mathcal{A}} \leq c^{\mathcal{A}} (Y^{\mathcal{A}})^2 d \log\left(1 + \frac{TX^2}{\lambda d}\right) + \frac{\lambda (Y^{\mathcal{A}})^2 T}{\lambda_{r_T}(G_T(0))},$$

*where $r_T = \text{rank}(G_T(0))$ and $\lambda_{r_T}$ is its smallest positive eigenvalue, $c^{\mathbf{r}} = 4$ and $c^{\mathbf{f}} = 1$.*

See proof in Appendix A. This bound suggests that to obtain a $\log(T)$ bound, $\lambda$ should not be chosen larger than about $\log(T)/T$, due to the second term, this is the *stringent regularization limitation*.

Choosing $\lambda = 1/T$ yields a first order regret of $2dY^2$ for the forward algorithm and $8dY^2$ for ridge regression (with clipping and prior knowledge of $Y$), which is at best twice the first order term from the lower bound. This suggests the presence of an optimality gap. Strikingly, Gaillard et al. [11] show that a non-regularized version of the forward algorithm achieves the optimal first order of $dY^2$. However, it also suffers from an important weakness: Indeed, the $(Y^{\mathcal{A}})^2 / \lambda_{r_T}(G_T(0))$ term in Corollary 2.3.1 is not uniformly bounded over feature sequences, but only on specific "well-behaved" features. In fact, double uniformity over features and observations is still an open question (see [11]).

## 2.3 Limitations in the stochastic setting

Now that we have recalled the main properties of the forward and ridge algorithms in the adversarial setup, we advocate for the need of a complementary analysis of the previous algorithms in the stochastic unbounded setting by unveiling some key limitations.

**Too unconstrained** The existing analysis being for a different setting, it naturally ignores crucial aspects of the stochastic setup. For instance, the quantity $Y$ is uninformative and may be substantial. Let us look at how the term $Y$ appears in the proofs of Azoury & Warmuth [3]. For ridge regression, the penultimate step to prove Theorem 2.1 writes:

$$L_T^{\mathbf{r}} - \min_\theta \left(L_T(\theta) + \lambda\|\theta\|_2^2\right) \leq \sum_{t=1}^T \underbrace{\left(x_t^\top \theta_{t-1} - y_t\right)^2 x_t^\top G_t^{-1} x_t}_{\text{first term}} \leq 4(Y^{\mathbf{r}})^2 \sum_{t=1}^T x_t^\top G_t^{-1} x_t. \quad (4)$$

In an adversarial setting, the "first term" cannot be controlled without assuming bounded predictions $|x_t^\top \theta_{t-1}| \leq Y^{\mathbf{r}}$, and doing so yields a bound $\left(x_t^\top \theta_{t-1} - y_t\right)^2 \leq 4(Y^{\mathbf{r}})^2$. In a stochastic setup however, we expect the term $\left(x_t^\top \theta_{t-1} - y_t\right)^2$ to reduce and stabilize around $\left(x_t^\top \theta_* - y_t\right)^2$, owing to the convergence properties of the estimate towards $\theta_*$.

For the forward algorithm, the final step in the proof of Theorem 2.2 writes:

$$L_T^{\mathbf{f}} - \min_\theta \left(L_T(\theta) + \lambda\|\theta\|_2^2\right) \leq \sum_{t=1}^T \underbrace{y_t^2 x_t^\top G_t^{-1} x_t}_{\text{first term}} - \sum_{t=1}^{T-1} \underbrace{\boldsymbol{x}_{t+1}^\top \boldsymbol{G}_t^{-1} \boldsymbol{x}_{t+1} \left(\boldsymbol{x}_{t+1}^\top \boldsymbol{\theta}_t\right)^2}_{\text{second term}}. \quad (5)$$

Then, the analysis uses that $|y_t| \leq Y^{\mathbf{f}}$ and disregards the negative contribution of the "second term". *Illustrative example:* Let us analyze these terms in an practice: consider $d = 5$, $\theta_* \in \mathbb{R}^5$, we sample 200 features uniformly in $[0,1]^5$ and Gaussian noises ($\sigma = 0.1$). Fig. 1 displays the instantaneous first regret term of both algorithms (with $\lambda = 1$) and the second regret term of the forward algorithm, averaged over 100 replicates. We remark that the first terms vanish quickly for ridge regression and are quite stable for the forward algorithm. On the other hand, they are essentially cancelled out by the second term.
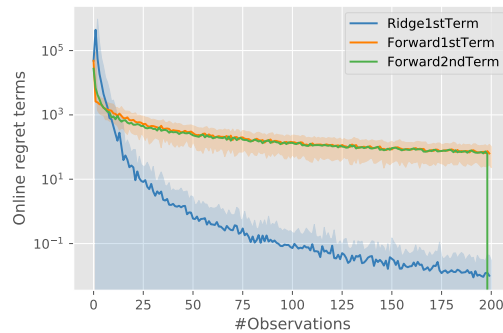


Figure 1: Online regret. $y$-axis is logarithmic.

5

Overall, the two strategies perform on par on this example. This suggests that Theorems 2.1 and 2.2 can be misleading in this stochastic setup: for ridge regression they introduce a conservative $4(Y^{\mathbf{r}})^2$ bound on $(x_t^\top(\theta_{t-1} - \theta_*))^2$, while in practice we observe that this term decreases rapidly to zero; for the forward algorithm, the bound ignores the effect of a negative term, which, as we see in Fig. 1, is essential to explain why this algorithm may outperform ridge regression.

**Time dependence:**    In the stochastic setup, it can be confusing to introduce $Y^{\mathcal{A}}$ because this hides a significant dependence on time. Indeed, for the forward algorithm $Y^{\mathbf{f}} = \max_{1 \le t \le T} |x_t^\top \theta_* + \epsilon_t|$. Considering the tractable setting of Gaussian i.i.d noise with variance $\sigma^2$, by classical Sudakov minoration from [20], we deduce that there exists $C > 0$ such that:

$$\forall T \ge 1 : \mathbb{E}[Y^{\mathbf{f}}] \ge \mathbb{E}\left[\max_{1 \le t \le T} \epsilon_t\right] - X\|\theta_*\|_2 \ge \sigma C \sqrt{2\log(T)} - X\|\theta_*\|_2.$$

Since $(Y^{\mathbf{f}})^2$ appears in the previous performance bounds, this suggests that $Y^{\mathbf{f}}$ actually increases the order of the regret bound to $\log(T)^2$ in this setting.

By focusing on the *unbounded stochastic* scenario, we hope in this paper to shed novel light on the practical performance of these strategies and better explain these phenomena.

## 3    High probability bounds

In this section, we analyze online ridge regression and the forward algorithm in the *stochastic* setting. We present our results in terms of the following intuitive regret definition:

$$\bar{R}_T^{\mathcal{A}} = L_T^{\mathcal{A}} - L_T(\theta_*). \tag{6}$$

This regret directly compares the cumulative loss of the learner to the cumulative loss of the oracle knowing the true parameter $\theta_*$. This contrasts with the online-to-batch conversion result that compares the loss of the learner to the value of a batch regularized optimization problem. Since we are in a stochastic setup, we further state results in high probability. More precisely, we state Theorems 3.1,3.2, 3.3 below holding with high probability uniformly over all $T$, and not simply for each $T$. As a first step, we prove that for $T$ great enough, we can choose this definition instead of $R_T$ defined in Eq. 1 without altering the bounds.

**Theorem 3.1.** *(Regret equivalence) In the stochastic setting with sub-Gaussian noise, for all $\delta > 0$ with probability at least $1 - \delta$, for all $T > 0$, $(x_t)_{1 \le t \le T} \in \mathbb{R}^d$ such that $\|x_t\|_2 \le X, |G_T(0)| > 0$*

$$L_T(\theta_*) - \min_{\theta \in \mathbb{R}^d} L_T(\theta) = o\left(\log(T)^2\right),$$

*in particular, it comes*

$$R_T^{\mathcal{A}} = \bar{R}_T^{\mathcal{A}} + o\left(\log(T)^2\right)$$

We detail the proof in Appendix B. Theorem 3.1 justifies choosing $\bar{R}_T$ to provide first order guarantees. Indeed, in the following sections we prove high probability upper bounds of order $O(\log(T)^2)$.

### 3.1    Online ridge regression

We start our results by stating a new high probability regret bound for online ridge regression.

**Theorem 3.2.** *In the stochastic setting with sub-Gaussian noise, for all $\delta > 0$ with probability at least $1 - \delta$, for all $T \ge 0$:*

$$\bar{R}_T^r \le \frac{2d\sigma^2 X^2}{\lambda \log(1 + X^2/\lambda)} \log\left(1 + TX^2/\lambda d\right) \log\left(\frac{(1 + TX^2/\lambda d)^{d/2}}{\delta/2}\right) + o(\log(T)^2),$$

*where $X = \max_{1 \le t \le T} \|x_t\|_2$.*

We refer to Appendix C for the proof and for the full regret expression. This result is interesting because the *ranges* of both predictions and observations do not appear, hence predictions clipping and/or a prior knowledge assumption on $Y^{\mathbf{f}}$ are not required. On the other hand, a factor $1/\lambda$ appears in the worst case of a singular design matrix. This seems to be the price for no longer assuming bounded predictions. Another notable improvement is that this bound no longer involves $\lambda\|\theta\|_2^2$ terms. In particular, it is uniform over bounded sequences of observations.

**Remark 2.** *(Regularization in ridge) Note that the bound holds with high probability, uniformly over $T$, and not only for each individual time horizon. In the proof of this result, $1/\lambda$ emerges from bounding $\lambda_{min}(G_t(0))$ in the worst case. When the collected features ensure the design matrix $G_t(0)$ is invertible, $1/\lambda$ virtually disappears. We highlight this experimentally in Section 3.4.*

## 3.2 The forward algorithm

We now derive a high probability regret bound for the forward algorithm using similar techniques.

**Theorem 3.3.** *Assuming sub-Gaussian noise, with probability at least $1 - \delta$, for all $T \geq 0$:*

$$\bar{R}_T^f \leq 2d\sigma^2 \log\left(1 + TX^2/\lambda d\right) \log\left(\frac{(1 + TX^2/\lambda d)^{d/2}}{\delta/2}\right) + o(\log(T)^2),$$

*where $X = \max_{1 \leq t \leq T} \|x_t\|_2$ and the $o(\log(T)^2)$ depends on $\lambda$ (see Appendix D).*

*Proof.* **(sketch)** We refer to Appendix D for the full proof and outline here the main steps leading to this bound. First, the instantaneous regret writes $\bar{r}_t = \ell_t(\theta_{t-1}) - \ell_t(\theta_*) = \left((\theta_{t-1} - \theta_*)^\top x_t\right)^2 + 2\epsilon_t(\theta_{t-1} - \theta_*)^\top x_t$. To control the term involving the noise $\epsilon_t$, we derive ad-hoc self-normalized tail inequalities in the same style of Theorem 1 in [1], which are of independent interest. Such an adaptation is required due to the forward algorithm. We then focus on controlling the first term. We construct confidence intervals for $\theta_*$ and resort to standard technical tools. $\square$

Theorem. 3.3 exhibits a better bound than Theorem. 3.2. In fact, the coefficient of the first order term for the forward algorithm only depends on the dimensionality and the noise variance, whilst for ridge regression, it also depends on the features' scale and on the regularization parameter $\lambda$.

**Remark 3.** *(Unrestrained regularization) Compared to existing results, this analysis lifts the "stringent regularization" that requires $\lambda = 1/T$ or data-dependent regularization (cf. [16]) to obtain uniform bounds. Therefore, Theorems 3.2 and 3.3 are not a mere consequence of bounding $Y^2$ with high probability in previous deterministic theorems. For completeness, we also derive a high probability regret bound for a non-regularized version of the forward algorithm in Appendix E; this algorithm was proven to be asymptotically first order minimax optimal in the adversarial bounded setting [11].*

## 3.3 Tightness of the bounds

Here we clarify the impact of a tighter confidence width for regularized least squares that was proved concurrently with the writing of this paper. First we state the result then we discuss its implications.

**Theorem 3.4.** *(Theorem 1 of Tirinzoni et al. [21]) Let $\delta \in (0, 1), n \geq 3$, and $\widehat{\theta}_t$ be a regularized least-square estimator obtained using $t \in [n]$ samples collected using an arbitrary bandit strategy $\pi := \{\pi_t\}_{t \geq 1}$. Then,*

$$\mathbb{P}\left\{\exists t \in [n] : \left\|\widehat{\theta}_t - \theta_*\right\|_{\bar{V}_t} \geq \sqrt{c_{n,\delta}}\right\} \leq \delta$$

*where $c_{n,\delta}$ is of order $\mathcal{O}(\log(1/\delta) + d\log\log n)$.*

This has important implications for Theorem 3.1, Theorem 3.2 and Theorem 3.3: in short it re-scales their regret upper-bounds from $R_T = O\left((d\sigma)^2 \log(T)^2\right)$ to $R_T = O\left(d\sigma^2 \log(T) \log\log(T)\right)$.
The first order $(d\sigma)^2 \log(T)^2$ in our results is the product of **1)** $d\log T$ from the elliptical lemma, for bounding the sum of feature norms and **2)** $\sigma^2 \log(T^d/\delta)$ the confidence ellipsoid width in the estimation of the regression parameter. It is the second term that is altered following the new result from Tirinzoni et al. [21]. These tighter confidence intervals change the upper bounds to $O(d\sigma^2 \log(T) \log\log(T))$. The latter matches the popular lower bounds in excess risk literature (see *e.g.* Theorem 1 in Mourtada [17]) up to sub-logarithmic terms suggesting *the optimality of the forward algorithm in the stochastic setting.*

## 3.4 Experiment

We provide experimental evidence supporting the fact that our novel high probability analysis better reflects the influence of regularization than results its adversarial counterpart.
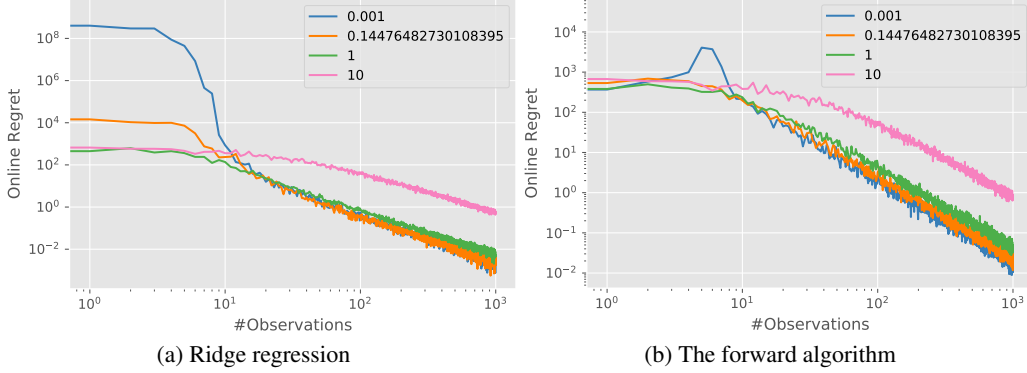


(a) Ridge regression

(b) The forward algorithm

Figure 2: Online regret's (Instantaneous loss difference) dependence on $\lambda$. All axes are logarithmic. Lines are averages over 100 repetitions and shaded areas represent one standard deviation.

In Figures 2a and 2b we observe the effect of regularization on the performance of ridge and forward regressions in a 5-dimensional regression setting, we vary $\lambda \in \{1/T, 1/\log(T), 1, 10\}$, sample a zero mean Gaussian noise with $\sigma = 0.1$ and draw features uniformly from the unit ball. The results clearly highlight the robustness of the forward algorithm to $\lambda$, contrarily to ridge. In particular, for ridge regression, we observe the exact dependence on $\lambda$ described by Theorem 3.2 in the first rounds of learning; as explained in Remark 2, once the collected features are enough for the design matrix $G_t(0)$ to become non-singular, the $1/\lambda$ virtually disappears from the first order regret bound and is replaced by the smallest eigenvalue of $G_t(0)$, making the regret significantly more stable.

## 4 Application: linear bandits

The proposed analysis of forward regression in the stochastic setting suggests that using it could be useful for revisiting several popular setups that include linear function approximation. We apply this change for stochastic linear bandits hereafter and derive the novel regret bound obtained when using forward regression instead of the standard ridge regression.

Consider the setting of *stochastic linear bandits*, where at round $t$ the reward of an action $x_t$ (from the action space $\mathcal{X} \subset \mathbb{R}^d$) is $y_t = \langle x_t, \theta_* \rangle + \epsilon_t$, where $\theta_* \in \mathbb{R}^d$ is an unknown parameter and $\epsilon_t$ is, conditionally on the past, a $\sigma$-sub-Gaussian noise. An upper bound $S$ on the unknown parameter's norm is provided: $\|\theta_*\|_2 \leq S$. The (pseudo) regret in this setting is defined:

$$R_T = \sum_{t=1}^{T} \langle x_t^*, \theta_* \rangle - \sum_{t=1}^{T} \langle x_t, \theta_* \rangle = \sum_{t=1}^{T} \langle x_t^* - x_t, \theta_* \rangle, \tag{7}$$

where $x_t^* = \arg\max_{x \in \mathcal{X}} \langle x, \theta_* \rangle$. Traditionally, the following additional assumption is made.

**Assumption 1.** *for all $x_t \in \mathcal{X}$    $\langle x_t, \theta_* \rangle \in [-1, 1]$.*

The "optimism in the face of uncertainty linear bandit" (OFUL) algorithm was introduced in [1]. OFUL resorts to ridge regression, constructs a confidence ellipsoid for the parameter estimate, and chooses the action that maximizes the upper-confidence bound on the reward. Under Assumption 1, [1] prove that the cumulative regret of OFUL satisfies, for $\delta > 0$ with probability at least $1 - \delta$, $\forall T > 0$

$$R_T^r \leq 4\sqrt{Td\log(\lambda + TX^2/d)}\left(\lambda^{1/2}S + \sigma\sqrt{2\log(1/\delta) + d\log(1 + TX^2/(\lambda d))}\right),$$

where $X = \max_{1 \leq t \leq T} \|x_t\|_2$.

**Forward variant [Algorithm 3]:** In a second phase, we propose the variant OFUL$^f$ in which we replace ridge regression by the forward algorithm. What this means is that the parameter estimate is a function of actions:

$$\theta_t^f(x) = \arg\min_{\theta \in \mathbb{R}^d} \sum_{s=1}^{t} (y_s - \langle x_s, \theta \rangle)^2 + \lambda \|\theta\|_2^2 + \langle x, \theta \rangle^2.$$

This fits perfectly because the new action can be chosen. Implementation details are in Algorithm 3.

---

**Algorithm 3:** OFUL$^f$ algorithm

---

Given $\lambda, \delta, S > 0$

**for** $t = 1, \ldots, T$ **do**

$\quad x_t = \arg\max_{x \in \mathcal{X}} \langle x, \theta_t^f(x) \rangle + \|x\|_{G_{t-1,x}^{-1}} (\sqrt{\lambda} + \|x\|_2)S + \sigma \sqrt{2 \log \left( \frac{(1 + tX_t^2(x)/\lambda d)^{d/2}}{\delta} \right)}$,

$\quad$ where $X_t(x) = \max\{\|x\|_2, \max_{1 \leq s \leq t-1} \|x_s\|_2\}$, $G_{t-1,x} = G_{t-1} + xx^\top$ and

$\quad \theta_t^f(x) = \arg\min_{\theta \in \mathbb{R}^d} \sum_{s=1}^{t-1} (y_s - \langle x_s, \theta \rangle)^2 + \lambda \|\theta\|_2^2 + \langle x, \theta \rangle^2$

$\quad$ play $x_t$ and observe $y_t$.

**end**

---

Note that OFUL$^f$ only requires an upper bound $S$ on $\|\theta_*\|_2$. We prove that OFUL$^f$ enjoys the same regret bound as OFUL and doesn't require Assumption 1. In stark contrast, we cannot show a similar bound for the standard OFUL without said assumption, it actually suffers a $\lambda$-dependent scaling factor in this case.

**Theorem 4.1.** *(Bandits with unbounded rewards) Without Assumption 1, for all $\delta > 0$, OFUL$^r$ achieves with probability at least $1 - \delta$, for all $T \geq 1$,*

$$R_T^r \leq 4 \sqrt{\frac{\boldsymbol{X}^2}{\boldsymbol{\lambda} \log(1 + \boldsymbol{X}^2/\boldsymbol{\lambda})} Td \log(\lambda + TX^2/d)} \left( \lambda^{1/2}S + \sigma\sqrt{2\log(1/\delta) + d\log(1 + TX^2/(\lambda d))} \right),$$

*also, we show that for all $\delta > 0$, OFUL$^f$ achieves with probability at least $1 - \delta$, for all $T \geq 1$:*

$$R_T^f \leq 4\sqrt{Td\log(\lambda + TX^2/d)} \left( (\lambda^{1/2} + X)S + \sigma\sqrt{2\log(1/\delta) + d\log(1 + TX^2/(\lambda d))} \right).$$

**Remark 4.** *we can drop the dependence on $X$ and $S$ by bounding the second term in the index of OFUL and OFUL$^f$ (see line 285) by $XS(1 + X/\sqrt{\lambda})$ and then dropping this -constant- term at the expense of a looser index. Therefore, knowing the bounds $x$ and $S$ is not crucial. Furthermore, while we choose to adopt the pseudo-regret definition like in [1], we could also derive similar bounds for the regret involving rewards $y_t = \langle x_t, \theta_* \rangle$ instead of their expected value, $(y_t)_{t \geq 1}$ are unbounded.*

**Experiment** We provide experimental evidence that the OFUL$^f$ variant improves OFUL for linear bandits; we find that it is generally as good as the standard OFUL$^r$, and in some cases it can prove to be significantly more robust to aberrant regularization parameters. We consider a 100-dimensional linear bandit with 10 arms, the parameter vector is drawn from the unit ball, actions are such that $\|x_t\| \leq 200$. Noise $\epsilon_t \overset{\mathcal{L}}{=} \mathcal{N}(0, 10^{-1})$, $\lambda = 10^{-5}$, $\delta = 10^{-3}$. In Fig. 3, lines are average regret over 100 repetitions and shaded areas cover the region between dashed-lines that are the first and third quartiles.



Figure 3: Cumulative regret. $y$-axis is logarithmic.

We observe that -as predicted by Theorem 4.1: OFUL$^f$ is particularly robust and choosing $\lambda = 1/T$ incurs substantial regret for OFUL. Because of this phenomena, and for the same observations in the online stochastic regression setting, we advocate for the use of the forward algorithm instead of ridge regression whenever possible, to take advantage of its increased robustness to $\lambda$.
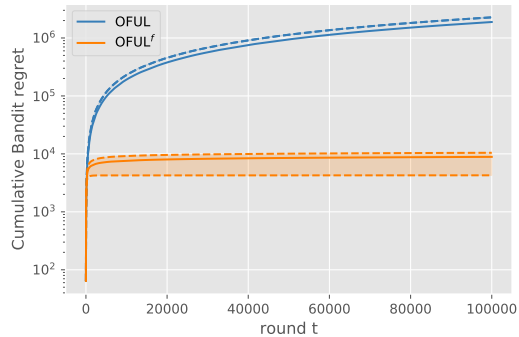
**Remark 5.** *Regarding the choice $\lambda = 1/T$: we use this specific regularization for two reasons: 1) to demonstrate the benefits of our stochastic analysis, since previous deterministic bounds suggest this $\lambda$ is best, 2) to showcase the increased robustness of OFUL$^{\mathsf{f}}$ compared to OFUL. In fact, more often than not, OFUL performs as good as OFUL, except when $\lambda$ is small or $X$ is large.*

## 5   Conclusion

We revisited the analysis of online linear regression algorithms in the setup of stochastic, possibly unbounded observations. We proved high probability regret bounds for three popular online regression algorithms (*cf.* Theorems 3.2, 3.3 and E.1). These bounds provide novel understanding of online regression. In particular, Theorem 3.2 seems to be the first regret bound for ridge regression that does not require bounded predictions or prior knowledge of a bound on observations. Our novel bounds seem to correctly capture the nature of dependence with regularization, as indicated by Fig. 2. Moreover, a new results from Tirinzoni et al. [21] can be incorporated in the proof mechanism to bring the high probability upper bounds to $O(d\sigma^2 \log(T) \log\log(T))$, which matches the optimal achievable bounds from the excess risk literature up to sub-logarithmic factors.

Furthermore, we argue that replacing ridge regression by the forward algorithm whenever possible in algorithms that require linear approximations can be beneficial, we depict this in a case study involving linear bandits: First from a theoretical standpoint our results show that the OFUL$^{\mathsf{f}}$ algorithm enjoys the classic first order regret bound while dropping Assumption 1; Second, we find that empirically, implementing OFUL with the forward algorithm makes the algorithm significantly more robust to extreme values of regularization, which is of practical interest.

More broadly, we believe that the improvement resulting from replacing ridge regression with the forward algorithm could be extended to several other settings: For instance, we also provide a similar analysis for non-stationary linear bandits in Appendix F; Graph bandits are of interest as well: they consider linear function approximations using ridge regression, and make Assumption 1, see for example Theorem 1 of [23]; Meta-learning with linear bandits can also be enhanced using forward regression: see for example Lemma 1 and consequent results in [5].

**Societal impact**   While our findings are purely theoretical, they can be taken advantage of for activities with negative societal impact. For instance, bandit algorithms are mainly used nowadays for advertising which can sometimes be linked with invasion of privacy. We insist however that linear regression in its generality is a very valuable technique, its use is ubiquitous is scientific domains, and improving our understanding of it even slightly is beneficial.

## Acknowledgments and Disclosure of Funding

## References

[1] Abbasi-Yadkori, Y., Pál, D., and Szepesvári, C. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, pp. 2312–2320, 2011.

[2] Abbasi-Yadkori, Y., Pal, D., and Szepesvari, C. Online-to-confidence-set conversions and application to sparse stochastic bandits. In *Artificial Intelligence and Statistics*, pp. 1–9. PMLR, 2012.

[3] Azoury, K. S. and Warmuth, M. K. Relative loss bounds for on-line density estimation with the exponential family of distributions. *Machine Learning*, 43(3):211–246, 2001.

[4] Bartlett, P. L., Koolen, W. M., Malek, A., Takimoto, E., and Warmuth, M. K. Minimax fixed-design linear regression. In *Conference on Learning Theory*, pp. 226–239, 2015.

[5] Cella, L., Lazaric, A., and Pontil, M. Meta-learning with stochastic linear bandits. In *International Conference on Machine Learning*, pp. 1360–1370. PMLR, 2020.

[6] Cesa-Bianchi, N. and Lugosi, G. *Prediction, learning, and games*. Cambridge university press, 2006.

[7] Cesa-Bianchi, N., Long, P. M., and Warmuth, M. K. Worst-case quadratic loss bounds for prediction using linear functions and gradient descent. *IEEE Transactions on Neural Networks*, 7(3):604–619, 1996.

[8] Cortes, C. and Mohri, M. On transductive regression. In *Advances in Neural Information Processing Systems*, pp. 305–312, 2007.

[9] Foster, D. P. Prediction in the worst case. *The Annals of Statistics*, pp. 1084–1090, 1991.

[10] Freedman, D. A. On tail probabilities for martingales. *the Annals of Probability*, pp. 100–118, 1975.

[11] Gaillard, P., Gerchinovitz, S., Huard, M., and Stoltz, G. Uniform regret bounds over $\mathbb{R}^d$ for the sequential linear regression problem with the square loss. In *Algorithmic Learning Theory*, pp. 404–432, 2019.

[12] Kamath, G. Bounds on the expectation of the maximum of samples from a gaussian. *URL http://www. gautamkamath. com/writings/gaussian max. pdf*, 2015.

[13] Kivinen, J. and Warmuth, M. K. Exponentiated gradient versus gradient descent for linear predictors. *information and computation*, 132(1):1–63, 1997.

[14] Littlestone, N., Long, P. M., and Warmuth, M. K. On-line learning of linear functions. In *Proceedings of the twenty-third annual ACM symposium on Theory of computing*, pp. 465–475, 1991.

[15] Maillard, O.-A. Self-normalization techniques for streaming confident regression. 2016.

[16] Malek, A. and Bartlett, P. L. Horizon-independent minimax linear regression. In *Advances in Neural Information Processing Systems*, pp. 5259–5268, 2018.

[17] Mourtada, J. Exact minimax risk for linear least squares, and the lower tail of sample covariance matrices. *arXiv preprint arXiv:1912.10754*, 2019.

[18] Orabona, F., Cesa-Bianchi, N., and Gentile, C. Beyond logarithmic bounds in online learning. In *Artificial intelligence and statistics*, pp. 823–831. PMLR, 2012.

[19] Russac, Y., Vernade, C., and Cappé, O. Weighted linear bandits for non-stationary environments. *arXiv preprint arXiv:1909.09146*, 2019.

[20] Sudakov, V. N. Gaussian measures, cauchy measures and $\varepsilon$-entropy. In *Soviet Math. Dokl*, volume 10, pp. 310–313, 1969.

[21] Tirinzoni, A., Pirotta, M., Restelli, M., and Lazaric, A. An asymptotically optimal primal-dual incremental algorithm for contextual linear bandits. *arXiv preprint arXiv:2010.12247*, 2020.

[22] Tripuraneni, N. and Mackey, L. Single point transductive prediction. *arXiv*, pp. arXiv–1908, 2019.

[23] Valko, M., Munos, R., Kveton, B., and Kocák, T. Spectral bandits for smooth graph functions. In *International Conference on Machine Learning*, pp. 46–54. PMLR, 2014.

[24] Vovk, V. Competitive on-line statistics. *International Statistical Review*, 69(2):213–248, 2001.

## A  Deterministic regret upper bound

In this section, we prove Corollary 2.3.1 in which we provide the explicit regret bounds for online ridge regression and forward regression in the adversarial case. First, recall the following result.

**Theorem.** *(Theorem 11.8 of Cesa-Bianchi & Lugosi [6]) For all $T \geq 1, (x_t)_{1 \leq t \leq T} \in \mathbb{R}^d$, $(y_t)_{1 \leq t \leq T} \in [-Y, Y]$ such that $\|x_t\|_2 \leq X$,*

$$\text{for } \mathcal{A} \in \{r, f\} \qquad \bar{R}_T^{\mathcal{A}} \leq c^{\mathcal{A}} (Y^{\mathcal{A}})^2 d \ln \left(1 + \frac{TX^2}{\lambda d}\right) + \lambda \|\theta_T\|_2^2,$$

*where $c^r = 4, c^f = 1, Y^r = \max\{Y, \max_{1 \leq t \leq T} |x_t^\top \theta_{t-1}^r|\}, Y^f = Y$ and $\theta_T = \arg\min_\theta L_T(\theta)$.*

We can now derive the explicit regret bound we seek by bounding the norm of the parameter $\theta_T$.

**Corollary.** *(Corollary 2.3.1) For all $T \geq 1, (x_t)_{1 \leq t \leq T} \in \mathbb{R}^d, (y_t)_{1 \leq t \leq T} \in [-Y, Y]$ such that $\|x_t\|_2 \leq X$,*

$$\text{for } \mathcal{A} \in \{r, f\} \qquad \bar{R}_T^{\mathcal{A}} \leq c^{\mathcal{A}} (Y^{\mathcal{A}})^2 d \ln \left(1 + \frac{TX^2}{\lambda d}\right) + \frac{\lambda (Y^{\mathcal{A}})^2 T}{\lambda_{r_T}(G_T(0))}$$

*where $r_T = \text{rank}(G_t(0))$ and $\lambda_{r_T}$ is its smallest positive eigenvalue, $c^r = 4, c^f = 1, Y^r = \max\{Y, \max_{1 \leq t \leq T} |x_t^\top \theta_{t-1}^r|\}, Y^f = Y$.*

*Proof.* Consider (w.l.o.g) ridge regression, denote $X_T$ the design matrix and $y_T$ the labels, then:

$$\|\theta_T\|_2 = \left\|G_T(0)^\dagger \mathbf{b}_T\right\|_2 = \sqrt{y_T^\top X_T^\top G_T(0)^\dagger G_T(0)^\dagger X_T y_T} \leq \sqrt{\frac{y_T^\top X_T^\top G_T(0)^\dagger X_T y_T}{\lambda_{r_T}(G_T)}}$$

$$\leq Y^r \sqrt{\frac{T}{\lambda_{r_T}(G_T)}},$$

where $G_T(0)^\dagger$ is the pseudo-inverse of $G_T(0)$, the last inequality is because $X_T^\top G_T(0)^\dagger X_T$ is an orthogonal projection on $\text{Im}(X^\top)$. Injecting in the previous theorem finishes the proof, these bounds hold for arbitrary bounded sequences. The proof for the forward algorithm proceeds in the same way by replacing $G_T$ by $G_{T+1}$ and $Y^r$ by $Y^f$. $\qquad \square$

## B  Regret definition

In this section, we prove that with high probability, $\bar{R}_T$ and $R_T$ yield the same first order high probability bounds for online regression algorithms.

**Theorem.** *(Regret equivalence) For all $\delta > 0$, with probability at least $1 - \delta$, for $T > 0$ such that $\sum_{s=1}^t x_s x_s^\top$ is non-singular:*
$$R_T = \bar{R}_T + o(\log(T)^2)$$

Note that this is enough to prove that $R_T$ and $\bar{R}_T$ are equal in first order because the upper bound on $\bar{R}_T$ is of order $\log(T)^2$.

Denote $\forall T \geq 1 : \theta_T = \arg\min_{\theta \in \mathbb{R}^d} L_T(\theta)$, then:

$$R_T - \bar{R}_T = L_T(\theta_*) - L_T(\theta_T) = 2 \sum_{t=1}^T \epsilon_t (\theta_T - \theta_*)^\top x_t - \sum_{t=1}^T \left((\theta_T - \theta_*)^\top x_t\right)^2. \qquad (8)$$

Denote $S_T = \sum_{t=1}^T \epsilon_t (\theta_T - \theta_*)^\top x_t, A_T = \sum_{t=1}^T \left((\theta_T - \theta_*)^\top x_t\right)^2$, we prove that $S_T = o(A_T)$.

**Lemma B.1.** *(Tail inequality) For all $\delta > 0, \sigma' > 0$, with probability at least $1 - \delta$, for all $T > 0$:*

$$|S_T| \leq \sqrt{2(A_T + 1/\sigma'^2) \log\left(\frac{\sqrt{\sigma'^2 A_T + 1}}{\delta}\right)}$$

*Proof.* We use the method of mixtures, denote

$$M_t^\lambda = \exp\left(\lambda \epsilon_t (\theta_T - \theta_*)^\top x_t - \frac{\lambda^2}{2}\left((\theta_T - \theta_*)^\top x_t\right)^2\right).$$

Without loss of generality, we can assume that $(\epsilon_s)_{s\geq 1}$ is 1-sub-Gaussian (this can be achieved by scaling features appropriately), then $\mathbb{E}[M_t^\lambda] \leq 1$.

Let $\Lambda \sim \mathbb{N}(0, \sigma'^2)$ be a Gaussian random variable and define $M_t = \mathbb{E}[M_t^\Lambda | F^\infty]$. We have $\mathbb{E}[M_t] = \mathbb{E}[\mathbb{E}[M_t^\Lambda | \Lambda]] \leq 1$. By making explicit $M_t$ and using Markov's inequality we get that for any stopping time $\tau$, for all $\delta > 0$, with probability at least $1 - \delta$:

$$\frac{|S_\tau|^2}{1/\sigma'^2 + A_\tau} \leq 2\sigma^2 \log\left(\frac{\sqrt{1 + \sigma'^2 A_\tau}}{\delta}\right).$$

We conclude using the same stopping time construction in Proof. C. $\qquad\square$

From Lemma. B.1 and equation. (8) we get that for all $\sigma', \delta > 0$ with probability at least $1 - \delta$:

$$
\begin{aligned}
R_T - \bar{R}_T &\leq \sqrt{2(A_T + 1/\sigma'^2)\log\left(\frac{\sqrt{\sigma'^2 A_T + 1}}{\delta}\right)} - A_T \\
&\leq \sqrt{(A_T + 1/\sigma'^2)\left(\log(\sigma'^2 A_T + 1) + 2\log(1/\delta)\right)} - A_T \\
&\leq \sqrt{A_T + 1/\sigma'^2}\left(\sqrt{\log(\sigma'^2 A_T + 1)} + \sqrt{2\log(1/\delta)}\right) - A_T \\
&\leq \frac{1}{\sigma'^2} + \sqrt{2(A_T + 1/\sigma'^2)\log(1/\delta)}
\end{aligned}
\tag{9}
$$

The next step is to the use confidence intervals of Maillard [15] which hold once the design matrix is singular.

**Theorem.** *(Theorem 3.3 of [15]) (Ordinary Least-squares) Assume that $N$ is a stopping time adapted to the filtration of the past. Then in the sub-Gaussian streaming regression model, for any $\delta > 0$, with probability at least $1 - \delta, \forall T \geq 1$ if $|G_T(0)| > 0$:*

$$\|\theta_* - \theta_T\|_{G_T(0)}^2 \leq 2(1 + \kappa)(1 + \alpha)\sigma^2 \log\frac{\kappa_d(e^2 \lambda_{max}(G_T))}{\delta}$$

*where $\kappa_d(x)$ is function of $\kappa$ and $\alpha$, $\kappa_d(x) = \frac{2}{3}\pi^2 \log(x/e)^2 \left[\frac{\log(x)}{2}\right]\left[(12(d+1)\sqrt{d})^d x^d + d\right]$ for $\kappa = \alpha = 1$.*

For bounded features $\|x\| \leq X$, we bound $\lambda_{max}(G_T(0)) \leq TX^2$. Denote $T_0 = \inf_{t\geq 1}\{|G_t| > 0\}$, and for $t \geq T_0 : \beta_t = 2(1 + \kappa)(1 + \alpha)\sigma^2 \log\frac{\kappa_d(e^2\lambda_{max}(G_T))}{\delta}$, then for all $\delta > 0$ with probability at least $1 - \delta$:

$$
\begin{aligned}
A_T = \sum_{t=1}^T \left((\theta_T - \theta_*)^\top x_t\right)^2 &\leq A_{T_0} + \sum_{t=T_0}^T \left((\theta_T - \theta_*)^\top x_t\right)^2 \\
&\leq A_{T_0} + \sum_{t=T_0}^T \beta_t \|x_t\|_{G_t(0)^{-1}}^2 \leq A_{T_0} + \beta_T \sum_{t=T_0}^T \|x_t\|_{G_t(0)^{-1}}^2
\end{aligned}
\tag{10}
$$

Then we bound the sum of features.

**Lemma B.2.** *(Technical inequality) For all sequences $\{x_t\}_t \in \mathbb{R}^d$ such that $\forall t, \|x_t\|_2 \leq X$, for all $\lambda \in \mathbb{R}_+, T_0, T \in \mathbb{N}$*

$$\sum_{t=T_0}^T \|x_t\|_{G_t^{-1}}^2 \leq d\log\left(1 + TX^2/\lambda_{min}(G_{T_0})d\right)$$

*where $G_t = G_t(\lambda)$.*

*Proof.* Using the Weinstein–Aronszajn identity: $\|x_t\|^2_{G_t^{-1}} = 1 - \frac{|G_{t-1}|}{|G_t|}$, and that $z - 1 \geq \log(z)$ leads to:

$$\sum_{t=T_0}^{T} \|x_t\|^2_{G_t^{-1}} \leq \sum_{t=1}^{T} - \log \frac{|G_{t-1}|}{|G_t|} = \log \left( \frac{|G_T|}{|G_{T_0}|} \right).$$

Since $\|x_t\|_2 \leq X$, using the AM-GM inequality:

$$\sum_{t=T_0}^{T} \log \left( 1 + \|x_t\|^2_{G_{t-1}^{-1}} \right) \leq d \log \left( 1 + TX^2/\lambda_{min}(G_{T_0})d \right).$$

$\square$

From equation (9), using $A_T \leq \sum_{t=1}^{T} \|\theta_T - \theta_*\|^2_{G_t} \|x_t\|^2_{G_t^{-1}} \leq \sum_{t=1}^{T} \|\theta_T - \theta_*\|^2_{G_T} \|x_t\|^2_{G_t^{-1}}$, then injecting Lemma B.2 with $\lambda = 0$, we find that for all $\delta > 0$, with probability at least $1 - \delta$:

$$A_T \leq A_{T_0} + \beta_T d \log \left( 1 + TX^2/\lambda_{min}(G_{T_0}(0))d \right)$$

Then injecting this last inequality in equation (8) gives, for all $\delta, \sigma' > 0$, with probability at least $1 - \delta$:

$$R_T - \bar{R}_T \leq \frac{1}{\sigma'^2} + \sigma' \sqrt{2 \log(1/\delta) \left( \beta_T d \log \left( 1 + TX^2/\lambda_{min}(G_{T_0})d \right) + 1 \right)}.$$

We also know -by definition- that $R_T \geq \bar{R}_T$. This concludes the proof for the equivalence of the two regret definitions.

## C   Ridge regression analysis

Here we prove a high probability time-uniform upper bound for online ridge regression. Let's recall the statement of the theorem that we prove.

**Theorem.** *(Theorem 3.2) For any $\delta > 0$, with probability at least $1 - \delta$, for all $T > 0$:*

$$\bar{R}_T^r \leq (d\sigma)^2 \frac{X^2/\lambda}{\log(1 + X^2/\lambda)} \log \left( \frac{1 + TX^2/\lambda d}{\delta/2} \right) \log \left( 1 + TX^2/\lambda d \right) + o(\log(T)^2)$$

See Eq. 14 for an explicit bound. In particular, the $o(\log(T)^2)$ term is $O(\log(T)^{3/2})$.

Let's write the instantaneous regret:

$$\bar{r}_t = \ell_t(\theta_{t-1}) - \ell_t(\theta_*) = \left( \theta_{t-1}^\top x_t - \theta_*^\top x_t \right)^2 + 2\epsilon_t(\theta_{t-1}^\top x_t - \theta_*^\top x_t) \tag{11}$$

The proof proceeds in three steps, that we detail hereafter and then we explain how to combine them for the final result.

First step: Confidence bound to control the concentration of $\theta_{t-1}$ around $\theta_*$. For this we use the confidence ellipsoid from Abbasi-Yadkori et al. [1].

**Theorem.** *(Confidence ellipsoid for ridge regression) For any $\delta > 0$, with probability at least $1 - \delta$, for all $t > 0$:*

$$\|\theta_t^r - \theta_*\|_{G_t} \leq \sqrt{\beta_t(\delta)} = \sigma \sqrt{d \log \left( \frac{1 + tX^2/\lambda d}{\delta} \right)} + \lambda^{1/2} S.$$

It comes, with probability at least $1 - \delta$:

$$(\theta_{t-1} - \theta_*)^\top x_t \leq \|x_t\|_{G_{t-1}^{-1}} \|\theta_{t-1} - \theta_*\|_{G_{t-1}} \leq \sqrt{\beta_{t-1}(\delta)} \|x_t\|_{G_{t-1}^{-1}}.$$

Then, since $\beta_t$ is non-decreasing:

$$L_t - L_t^* \leq \beta_{T-1} \sum_{t=1}^{T} \|x_t\|_{\eta_{t-1}}^2 + 2 \sum_{t=1}^{T} \epsilon_t (\theta_{t-1} - \theta_*)^\top x_t. \tag{12}$$

Second step: Next we bound the sum of feature norms. The main idea here is to use linear algebra techniques to obtain a telescopic sum.

Lemma B.2 doesn't apply here because we have $\|x_t\|_{G_{t-1}^{-1}}$ instead of $\|x_t\|_{G_t^{-1}}$. We derive a similar lemma for this sum of feature norms.

**Lemma C.1.** *(Technical inequality) For all sequences $\{x_t\}_t \in \mathbb{R}^d$ such that $\forall t, \|x_t\|_2 \leq X$, for all $\lambda \in \mathbb{R}_+, T \in \mathbb{N}$*

$$\sum_{t=1}^{T} \|x_t\|_{G_{t-1}^{-1}}^2 \leq \frac{X^2/\lambda}{\log(1 + X^2/\lambda)} d \log\left(1 + TX^2/\lambda d\right)$$

*Proof.* We use the Weinstein–Aronszajn identity: $\|x_t\|_{G_{t-1}^{-1}}^2 = \frac{|G_t|}{|G_{t-1}|} - 1$, which leads to:

$$\sum_{t=1}^{T} \log\left(1 + \|x_t\|_{G_{t-1}^{-1}}^2\right) = \log\left(\frac{G_T}{G_0}\right).$$

Then since $\|x_t\|_2 \leq X$ and using the AM-GM inequality:

$$\sum_{t=1}^{T} \log\left(1 + \|x_t\|_{G_{t-1}^{-1}}^2\right) \leq d \log\left(1 + TX^2/\lambda d\right).$$

This next part is what differs from Lemma B.2, using $\|x_t\|_{G_{t-1}^{-1}}^2 \leq \lambda_{\max}(G_{t-1}^{-1})\|x_t\|_2^2 \leq X^2/\lambda$ and the concavity of the function $\log$ we find:

$$\sum_{t=1}^{T} \|x_t\|_{G_{t-1}^{-1}}^2 \leq \sum_{t=1}^{T} \frac{X^2/\lambda}{\log(1 + X^2/\lambda)} \log\left(1 + \|x_t\|_{G_{t-1}^{-1}}^2\right).$$

The last inequality can also be proved by noting that $x \to x/\log(1+x)$ is non-decreasing which can be used to bound every feature norm. $\square$

Third step: To control the second term in the r.h.s of Eq. 11, we use Martingale inequalities similar to the ones used for the confidence intervals to derive a uniform high probability bound.

**Lemma C.2.** *(Tail inequality, see Corollary 8 of [2]) Define $S_t = \sum_{s=1}^{t} \epsilon_s (\theta_{s-1} - \theta_*)^\top x_s$ and let $(F_t)_{t \geq 0}$ be a filtration such that $x_t$ is $F_{t-1}$ measurable and $\epsilon_t$ is $F_t$ measurable. Then $S_t$ is a martingale with respect to $F_t$ and for any $\delta > 0, \sigma' > 0$, with probability at least $1 - \delta$, for all $t \geq 0$:*

$$|S_t| \leq \sigma \sqrt{2\left(1/\sigma'^2 + \sum_{s=1}^{t}\left((\theta_{t-1} - \theta_*)^\top x_t\right)^2\right) \log\left(\frac{\sqrt{1 + \sigma'^2 \sum_{s=1}^{t}\left((\theta_{t-1} - \theta_*)^\top x_t\right)^2}}{\delta}\right)}$$

*Proof.* The proof of this result follows the same line in the proof of Theorem 1 of Abbasi-Yadkori et al. [1], first we define for $\lambda \in \mathbb{R}^d, t > 0$ : $M_t^\lambda = \exp\left(\sum_{s=1}^{t}\left[\epsilon_s \lambda(\theta_{t-1} - \theta_*)^\top x_t - \lambda^2\left((\theta_{t-1} - \theta_*)^\top x_t\right)^2/2\right]\right)$.

Without loss of generality, we can assume that $(\epsilon_s)_{s \geq 1}$ is 1-sub-Gaussian (this can be achieved by scaling features). Let $\tau$ be a stopping time with respect to the filtration $\{F_t\}_{t=0}^{\infty}$. Then $M_\tau^\lambda$ is well-defined almost surely and

$$\mathbb{E}[M_\tau^\lambda] \leq 1.$$

15

Let $\Lambda \sim \mathbb{N}(0, \sigma'^2)$ be a Gaussian random variable and define $M_t = \mathbb{E}[M_t^\Lambda | F^\infty]$. We have $\mathbb{E}[M_t] = \mathbb{E}[\mathbb{E}[M_t^\Lambda | \Lambda]] \leq 1$. By expliciting $M_t$ and using Markov's inequality we get that for $\delta > 0$, with probability $1 - \delta$:

$$|S_\tau|^2 \leq \left( 1/\sigma'^2 + \sum_{t=1}^\tau \left( (\theta_{t-1} - \theta_*)^\top x_t \right)^2 \right) 2\sigma^2 \log \left( \frac{\sqrt{1 + \sigma'^2 \sum_{t=1}^\tau \left( (\theta_{t-1} - \theta_*)^\top x_t \right)^2}}{\delta} \right). \tag{13}$$

Next we use a stopping time construction from Freedman [10]: Define the bad event:

$$B_t(\delta) = \left\{ \omega \in \Omega : \frac{|S_t|^2}{1/\sigma'^2 + \sum_{s=1}^t \left( (\theta_{s-1} - \theta_*)^\top x_s \right)^2} > 2\sigma^2 \log \left( \frac{\sqrt{1 + \sigma'^2 \sum_{s=1}^t \left( (\theta_{s-1} - \theta_*)^\top x_s \right)^2}}{\delta} \right) \right\}$$

We are interested in bounding the probability that $\bigcup_{t > 0} B_t(\delta)$ happens. Define $\tau(\omega) = \min\{t \geq 0 : \omega \in B_t(\delta)\}$, with the convention that $\min \emptyset = \infty$. Then, $\tau$ is a stopping time. Further,

$$\bigcup_{t \geq 0} B_t(\delta) = \{\omega : \tau(\omega) < \infty\}$$

Thus, by Eq. 13:

$$\Pr \left[ \bigcup_{t \geq 0} B_t(\delta) \right] = \Pr[\tau < \infty] = \Pr[B_\tau(\delta), \tau < \infty] \leq \Pr[B_\tau(\delta)] \leq \delta$$

$\square$

This proves that the second term in Eq. 11 is a of order $\sim O(\log(T) \log \log T)$. In fact, with high probability $\sum_{s=1}^t \left( (\theta_{t-1} - \theta_*)^\top x_t \right)^2 = O(\log(T)^2)$ therefore, with high probability $S_T$ is of order $\sim O(\log(T) \log(\log T)/\delta)$. Consequently, with high probability, $S_T$ is second order.

Proof aggregation: By combining earlier results we find for any $\delta, \sigma' > 0$, with probability at least $1 - \delta$, for all $T \geq 0$:

$$\bar{R}_T^r \leq \left( \sigma \sqrt{d \log \left( \frac{1 + TX^2/\lambda d}{\delta/2} \right)} + \lambda^{1/2} S \right)^2 \frac{X^2/\lambda}{\log(1 + X^2/\lambda)} d \log \left( 1 + TX^2/\lambda d \right)$$

$$+ \sigma \sqrt{2 \left( 1/\sigma'^2 + \sum_{s=1}^t \left( (\theta_{t-1} - \theta_*)^\top x_t \right)^2 \right) \log \left( \frac{\sqrt{1 + \sigma'^2 \sum_{s=1}^t \left( (\theta_{t-1} - \theta_*)^\top x_t \right)^2}}{\delta/2} \right)}. \tag{14}$$

# D   Analysis of the forward algorithm

In this section we derive the high probability time-uniform regret bound for the forward algorithm. Let's recall the theorem.

**Theorem.** *(Theorem 3.3)For any $\delta > 0$, with probability at least $1 - \delta$, for all $T > 0$:*

$$\bar{R}_T^f \leq (d\sigma)^2 \log \left( \frac{1 + TX^2/\lambda d}{\delta/2} \right) \log \left( 1 + TX^2/\lambda d \right) + o(\log(T)^2)$$

See Eq. 15 for the explicit expression of this bound. The proof proceeds similarly to Appendix C: we need to bound the instantaneous regret.

$$\bar{r}_t = \ell_t(\theta_{t-1}) - \ell_t(\theta_*) = \left( \theta_{t-1}^\top x_t - \theta_*^\top x_t \right)^2 + 2\epsilon_t (\theta_{t-1}^\top x_t - \theta_*^\top x_t)$$

We proceed in three steps like before.

First step: We start by deriving a confidence ellipsoid for this new parameter estimate. This is a novel result.

**Theorem.** *(Confidence ellipsoid for the Forward algorithm) For any $\delta > 0$, with probability at least $1 - \delta$, for all $t > 0$:*

$$\|\theta_t - \theta_*\|_{G_t} \le \sqrt{\beta_t(\delta)} = \sigma\sqrt{d \log\left(\frac{1 + tX^2/\lambda d}{\delta}\right)} + (\lambda^{1/2} + X)S.$$

*Proof.* Denote $X_t = (x_1^\top, \dots, x_t^\top), \varepsilon_t = (\epsilon_1, \dots, \epsilon_t)^\top$. Using

$$\theta_t = G_{t+1}^{-1} X_t^\top (X\theta_* + \varepsilon_t) = G_{t+1}^{-1} X_t^\top \varepsilon_t + G_{t+1}^{-1}(X_t^\top X_t + \lambda I + x_{t+1}^\top x_{t+1})\theta_* - G_{t+1}^{-1}(\lambda I + x_{t+1}^\top x_{t+1})\theta_*$$
$$= G_{t+1}^{-1} X_t^\top \varepsilon_t + \theta_* - G_{t+1}^{-1}(\lambda I + x_{t+1}^\top x_{t+1})\theta_*,$$

we get

$$|x^\top \theta_t - x^\top \theta_*| = |x^\top G_{t+1}^{-1} X_t \varepsilon_t - x^\top G_{t+1}^{-1}(\lambda\theta_* + x_{t+1}x_{t+1}^\top\theta_*)|$$
$$\le \|x\|_{G_{t+1}^{-1}}\left(\|X_t^\top \varepsilon_t\|_{G_{t+1}^{-1}} + (\sqrt{\lambda} + X)\|\theta_*\|_2\right),$$

where in the last inequality we used Cauchy-Schwartz inequality and that by the Sherman-Morrison formula $x_{t+1}^\top G_{t+1}^{-1} x_{t+1} = \frac{x_{t+1}^\top G_t^{-1} x_{t+1}}{1 + x_{t+1}^\top G_t^{-1} x_{t+1}} \le 1$. We know that: $\|X_t^\top \varepsilon_t\|_{G_{t+1}^{-1}} \le \|X_t^\top \varepsilon_t\|_{G_t^{-1}}$ which allows us to use Theorem 1 from Abbasi-Yadkori et al. [1] that we recall just after this proof. We conclude by plugging $x = G_{t+1}(\theta_t - \theta_*)$. $\quad\square$

**Theorem.** *(Self-Normalized Bound for Vector-Valued Martingales). Let $\{F_t\}_{t=0}^\infty$ be a filtration. Let $\{\eta_t\}_{t=1}^\infty$ be a real-valued stochastic process such that $\eta_t$ is $F_t$ -measurable and $\eta_t$ is conditionally $R$ -sub-Gaussian for some $R \ge 0$ i.e.*

$$\forall \lambda \in \mathbb{R} \quad \mathbf{E}\left[e^{\lambda \eta_t} \mid F_{t-1}\right] \le \exp\left(\frac{\lambda^2 R^2}{2}\right)$$

*Let $\{X_t\}_{t=1}^\infty$ be an $\mathbb{R}^d$ -valued stochastic process such that $X_t$ is $F_{t-1}$ -measurable. Assume that $V$ is a $d \times d$ positive definite matrix. For any $t \ge 0$, define*

$$\bar{V}_t = V + \sum_{s=1}^t X_s X_s^\top \quad S_t = \sum_{s=1}^t \eta_s X_s.$$

*Then, for any $\delta > 0$, with probability at least $1 - \delta$, for all $t \ge 0$,*

$$\|S_t\|_{\bar{V}_t^{-1}}^2 \le 2R^2 \log\left(\frac{\det\left(\bar{V}_t\right)^{1/2} \det(V)^{-1/2}}{\delta}\right).$$

*Note that the deviation of the martingale $\|S_t\|_{\bar{V}_t^{-1}}^2$ is measured by the norm weighted by the matrix $\bar{V}_t^{-1}$ which is itself derived from the martingale, hence the name "self-normalized bound".*

For the first term, with probability at least $1 - \delta$ for all $t \ge 0$:

$$(\theta_{t-1} - \theta_*)^\top x_t \le \|x_t\|_{G_t^{-1}} \|\theta_{t-1} - \theta_*\|_{G_t}$$
$$\le \sqrt{\beta_{t-1}(\delta)}\|x_t\|_{G_t^{-1}} \le \sqrt{\beta_{T-1}(\delta)}\|x_t\|_{G_t^{-1}}.$$

Second step: We can use Lemma B.2 to bound the sum of feature norms. It comes

$$\sum_{t=1}^T (\theta_{t-1}^\top x_t - \theta_*^\top x_t)^2 \le \beta_T(\delta) d \log\left(1 + TX^2/\lambda d\right)$$

Third step: Again, we derive a high probability bound To control the second term in the r.h.s of (11).

**Lemma D.1.** *(Tail inequality) Define $S_t = \sum_{s=1}^t \epsilon_s(\theta_{s-1} - \theta_*)^\top x_s$ and let $(F_t)_{t\ge 0}$ be a filtration such that $x_t$ is $F_{t-1}$ measurable and $\epsilon_t$ is $F_t$ measurable. Then $S_t$ is a martingale with respect to $F_t$ and for any $\delta > 0, \sigma' > 0$, with probability at least $1 - \delta$, for all $t \ge 0$:*

$$|S_t| \le \sigma\sqrt{2\left(1/\sigma'^2 + \sum_{s=1}^t \left((\theta_{t-1} - \theta_*)^\top x_t\right)^2\right) \log\left(\frac{\sqrt{1 + \sigma'^2 \sum_{s=1}^t \left((\theta_{t-1} - \theta_*)^\top x_t\right)^2}}{\delta}\right)}$$

*Proof.* The proof of this result proceeds in the exact same way as for Lemma C.2. □

Proof aggregation: We combine previous results to finish the proof of the forward algorithm regret bound. For any $\delta, \sigma' > 0$, with probability at least $1 - \delta$, for all $T \geq 0$:

$$
\bar{R}_T^{\mathrm{f}} \leq \left( \sigma \sqrt{d \log \left( \frac{1 + TX^2/\lambda d}{\delta/2} \right)} + (\sqrt{\lambda} + X)S \right)^2 \frac{X^2/\lambda}{\log(1 + X^2/\lambda)} d \log \left( 1 + TX^2/\lambda d \right)
$$
$$
+ \sigma \sqrt{2 \left( 1/\sigma'^2 + \sum_{s=1}^t \left( (\theta_{t-1} - \theta_*)^\top x_t \right)^2 \right) \log \left( \frac{\sqrt{1 + \sigma'^2 \sum_{s=1}^t \left( (\theta_{t-1} - \theta_*)^\top x_t \right)^2}}{\delta} \right)}.
$$
(15)

# E    The unregularized-forward algorithm

For the sake of completeness, we propose a high probability bound on the regret of a non-regularized forward algorithm -studied in the adversarial bounded case in Gaillard et al. [11]- which achieves the optimal asymptotic first order deterministic minimax bound of $dY^2 \log(T)$. This algorithm is a simple yet elegant modification of forward regression, it avoids the exploding $\lambda \|\theta_T\|_2^2$ term by setting $\lambda = 0$. Consequently $\theta_t = G_{t+1}^\dagger b_t$, where $G_t^\dagger$ is the pseudo-inverse of $G_t$.

**Theorem E.1.** *(Regret of the unregularized forward) The unregularized forward regression achieves, for any $\delta > 0$, with probability at least $1 - \delta$ for all $T > 0$:*

$$
\bar{R}_T^{u-f} \leq 2(1 + \kappa)(1 + \alpha)\sigma^2 \log \left( \frac{\kappa_d(1 + TX^2/\gamma d)}{\delta/4} \right) \log \left( \frac{|G_T^\dagger|}{|G_{T_1}^\dagger|} \right)
$$
$$
+ 2\sigma^2 \log \left( \frac{4T_1}{\delta} \right) \left( d + \sum_{1 \leq t \leq T_1, t \in \mathcal{T}} \log \left( \frac{X^2}{\lambda_{r_t}(\sum_{s=1}^t x_t x_t^\top)} \right) \right),
$$

*where $\kappa, \alpha \in \mathbb{R}_+^*$ are peeling parameters (can be chosen), $\gamma = \min_{1 \leq t \leq T} \|x_t\|_2$, and $\kappa_d(x) \propto x^d$ up to logarithmic factors and depends on $\kappa$ and $\alpha$ (cf. Theorem 5.4 in Maillard [15]). $T_1 = \min \{t \geq 1, |G_t| > 0\}$ is, if it exists, the first time the design matrix is non-singular, otherwise $T_1 = T$, and $\mathcal{T}$ is the set of indices $t$ such that $rank(G_t) > rank(G_{t-1})$. The last term accounts for when the design matrix is singular, and is naturally unbounded (this was also the case in the adversarial case).*

Asymptotically, with probability at least $1 - \delta$ the first regret term is bounded as:

$$
\bar{R}_T^{u-f} \leq 2(1 + \kappa)(1 + \alpha) \log \left( \frac{C(\kappa, \alpha)(TX^2/\lambda d)^d}{\delta} \right) \log \left( (T - T_1)X^2/\lambda d \right),
$$

where $C(\kappa, \alpha)$ is a function of the peeling parameters.
We don't seek a more involved analysis to explicit this bound or improve on it, but we see that vaguely it leads to a bound similar to Theorems 3.2 and 3.3 provided that the term accounting for the singularity of the design matrix is controlled. The latter empowers the intuition that in the high probability analysis, the forward algorithm is *first order minimax optimal* even though concretely we cant be sure because we don't have access to uniform lower bounds.

*Proof.* The proof consists of two mains steps: the first is to use the following bound while the design matrix is singular:

**Theorem E.2.** *(Theorem 11 Gaillard et al. [11]) For all $T \geqslant 1$, for all sequences $x_1, \ldots, x_T \in \mathbb{R}^d$ and all $y_1, \ldots, y_T \in [-Y, Y]$, the unregularized forward algorithm achieves the regret bound*

$$
R_T(\mathbf{u}) \leq Y^2 \sum_{t=1}^T \mathbf{x}_t^\top \eta_t^\dagger \mathbf{x}_t \leqslant dY^2 \log T + dY^2 + Y^2 \sum_{t \in [1,T] \cap \mathcal{T}} \log \left( \frac{X^2}{\lambda_{r_t}(\sum_{s=1}^t x_s x_s^\top)} \right)
$$

*where $\forall M \in \mathcal{M}_d(\mathbb{R}), \lambda_1(M) \geq \ldots \geq \lambda_d$ are $M$'s eigenvalues and $r_t = \mathrm{rank}(\sum_{s=1}^t x_s x_s^\top))$ and where the set $\mathcal{T}$ contains $r_T$ rounds, given by the smallest $s \geqslant 1$ such that $\mathrm{x}_s$ is not mull, and all the $s \geqslant 2$ for which $\mathrm{rank}\,(\mathrm{G}_{s-1}) \neq \mathrm{rank}\,(\mathrm{G}_s)$.*

The second step is a bound when the design matrix is invertible, using Theorem B. Denote $T_1 = \inf_{t \geq 1}\{|G_t| > 0\}$, using Theorem E.2:

$$\bar{R}_{T_1} \leq Y^2 \left( d\log(T_1) + d + \sum_{1 \leq t \leq T_1, t \in \mathcal{T}} \log\left(\frac{X^2}{\lambda_{r_t}(G_t)}\right) \right)$$

From standard results on sub-Gaussian noise, we also know that $\mathbb{E}[\max_{1 \leq t \leq T} \epsilon_t] \leq \sigma\sqrt{2\log(T)}$ (see *e.g.* Kamath [12]), then using the transformation of Laplace along with Markov's inequality, $\forall \delta > 0\ \mathrm{P}\left(\forall T \geq 1, Y^2 \leq 2\sigma^2 \log(T/\delta)\right) \geq 1 - \delta$, hence with probability at least $1 - \delta$:

$$\bar{R}_{T_1} \leq 2d\sigma^2 \log\frac{T_1}{\delta}\log(T_1) + 2d\sigma^2 \log\frac{T_1}{\delta} + 2\sigma^2 \frac{\log T_1}{\delta} \sum_{1 \leq t \leq T_1, t \in \mathcal{T}} \log\left(\frac{X^2}{\lambda_{r_t}(\sum_{s=1}^t x_s x_s^\top)}\right). \tag{16}$$

And for $T > T_1$, we bound $R_T - R_{T_1}$ using the same methodology in Appendix C and Appendix D and using the confidence bounds above (*cf.* Theorem B). $\forall \delta > 0$, with probability at least $1 - \delta$:

$$\forall t > T_1 : \left(\theta_{t-1}^\top x_t - \theta_*^\top x_t\right)^2 \leq \sqrt{\beta_{t-1}(\delta)}\|x_t\|_{G_t^\dagger}$$

We use the tail inequality. (C.2) to get, $\forall \delta > 0$, with probability at least $1 - \delta, \forall T > 0$:

$$\bar{R}_T - \bar{R}_{T_1} \leq 2(1+\kappa)(1+\alpha)\sigma^2 \log\left(\frac{\kappa_d(1+TX^2/\lambda d)}{\delta/2}\right) \log\left(\frac{|G_T^\dagger|}{|G_{T_1}^\dagger|}\right) \tag{17}$$

From (16) and (17) we obtain for all $\delta > 0$, with probability at least $1 - \delta$:

$$\bar{R}_T \lesssim 2(1+\kappa)(1+\alpha)\sigma^2 \log\left(\frac{\kappa_d(1+TX^2/\lambda d)}{\delta/4}\right) \log\left(\frac{|G_T^\dagger|}{|G_{T_1}^\dagger|}\right)$$
$$+ 2\sigma^2 \frac{\log(T_1)}{\delta/4}\left(d + \sum_{1 \leq t \leq T_1, t \in \mathcal{T}} \log\left(\frac{X^2}{\lambda_{r_t}(\sum_{s=1}^t x_t x_t^\top)}\right)\right).$$

$\square$

# F  Applications

In this section, we provide technical details regarding the settings of stationary and non-stationary linear bandits.

## F.1  Linear bandits (Proof of Theorem 4.1)

We start by analyzing linear bandits in the stationary setting. Let us first see how OFUL$^{\mathrm{f}}$ behaves in the "unbounded rewards" scenario.

### F.1.1  OFUL with forward regression

Consider the same setting as that of Abbasi-Yadkori et al. [1], that we detailed in Section 4, we write the confidence interval $C_t(x)$ for the forward algorithm at the action $x$ as:

$$\left\{\theta \in \mathbb{R}^d : \|\theta_t^{\mathrm{f}} - \theta\|_{G_t + xx^\top} \leq \sqrt{\beta_t(x,\delta)} = (\sqrt{\lambda} + \|x\|_2)S + \sigma\sqrt{2\log\left(\frac{(1+tX^2/\lambda d)^{d/2}}{\delta}\right)}\right\}$$

which gives, for all $T \geq 0$ the regret (*cf.* Theorem 4.1):

$$R_T \leq 4\sqrt{Td\log(\lambda + TX^2/d)}\left(\lambda^{1/2}(S+X) + \sigma\sqrt{2\log(1/\delta) + d\log(1 + TX^2/(\lambda d))}\right)$$

this is equivalent to ridge in its first order, with better scaling and dependence on $\lambda$.

*Proof.* Lets decompose the instantaneous regret as follows:

$$r_t = \langle \theta_*, x_* \rangle - \langle \theta_*, x_t \rangle \leq \left\langle \tilde{\theta}_t, x_t \right\rangle - \langle \theta_*, x_t \rangle = \left\langle \tilde{\theta}_t - \theta_*, x_t \right\rangle$$

$$= \left\langle \widehat{\theta}_{t-1} - \theta_*, x_t \right\rangle + \left\langle \tilde{\theta}_t - \widehat{\theta}_{t-1}, x_t \right\rangle$$

$$= \left\| \widehat{\theta}_{t-1} - \theta_* \right\|_{(G_{t-1}+x_t x_t^\top)} \|X_t\|_{(G_{t-1}+x_t x_t^\top)^{-1}} + \left\| \tilde{\theta}_t - \widehat{\theta}_{t-1} \right\|_{(G_{t-1}+x_t x_t^\top)} \|x_t\|_{(G_{t-1}+x_t x_t^\top)^{-1}}$$

$$\leq 2\sqrt{\beta_{t-1}(x_t,\delta)} \|x_t\|_{(G_{t-1}+x_t x_t^\top)^{-1}}, \tag{18}$$

where $\widetilde{\theta}_t$ is the optimistic parameter estimate, *i.e.* the $\theta \in C_t(x_t)$ that maximizes the upper confidence bound on the reward of action $x_t$. The first inequality is since $\left(X_t, \widetilde{\theta}_t\right)$ is optimistic, and the last step holds by Cauchy-Schwarz. Using inequality (18) and the fact that $\sqrt{\beta_t(x,\delta)} \leq \sqrt{\beta_t(\delta)} = (\sqrt{\lambda} + X)S + \sigma\sqrt{2\log\left(\frac{(1+tX^2/\lambda d)^{d/2}}{\delta}\right)}$ we get that, with probability at least $1 - \delta$, for all $n \geq 0$

$$R_n \leq \sqrt{n\sum_{t=1}^n r_t^2} \leq \sqrt{8\beta_n(\delta)n\sum_{t=1}^n \|x_t\|_{(G_{t-1}+x_t x_t^\top)^{-1}}}$$

$$\leq 4\sqrt{nd\log(\lambda + nL/d)}\left((\lambda^{1/2} + X)S + \sigma\sqrt{2\log(1/\delta) + d\log(1 + nL/(\lambda d))}\right)$$

where the last step follow from Lemma B.2. $\qquad\square$

### F.1.2 OFUL with ridge regression

In this section, we derive a novel regret bound for online ridge regression, one that doesn't require the bounded rewards assumption (*cf.* Assumption 1).

**Theorem.** *(Bandits with unbounded rewards) Without Assumption 1, for all $\delta > 0$, OFUL$^r$ achieves with probability at least $1 - \delta$, for all $T \geq 1$,*

$$R_T^r \leq 4\sqrt{\frac{\boldsymbol{X}^2 Td\log(1 + TX^2/\lambda d)}{\boldsymbol{\lambda}\log(1 + \boldsymbol{X}^2/\boldsymbol{\lambda})}}\left(\lambda^{1/2}S + \sigma\sqrt{2\log(1/\delta) + d\log(1 + TX^2/(\lambda d))}\right),$$

*Proof.* The proof follows exactly like in Section F.1.1 except the last step (control of the norm of actions) that now proceeds using Lemma C.1. The first step is to use the confidence ellipsoid for the ridge regression parameter (see the second theorem in Section C or Theorem 2 of Abbasi-Yadkori et al. [1]). With probability at least $1 - \delta$, for all $t \geq 0, \theta_*$ lies in the set

$$C_t = \left\{ \theta \in \mathbb{R}^d : \|\theta_t^r - \theta\|_{G_t} \leq \sqrt{\beta_t(\delta)} = \sigma\sqrt{d\log\left(\frac{1+tX^2/\lambda d}{\delta}\right)} + \lambda^{1/2}S \right\}.$$

Then

$$r_t = \langle \theta_*, x_* \rangle - \langle \theta_*, x_t \rangle \leq \left\langle \tilde{\theta}_t, x_t \right\rangle - \langle \theta_*, x_t \rangle = \left\langle \tilde{\theta}_t - \theta_*, x_t \right\rangle$$

$$= \left\langle \widehat{\theta}_{t-1} - \theta_*, x_t \right\rangle + \left\langle \tilde{\theta}_t - \widehat{\theta}_{t-1}, x_t \right\rangle$$

$$= \left\| \widehat{\theta}_{t-1} - \theta_* \right\|_{G_{t-1}} \|X_t\|_{G_{t-1}^{-1}} + \left\| \tilde{\theta}_t - \widehat{\theta}_{t-1} \right\|_{G_{t-1}} \|x_t\|_{G_{t-1}^{-1}}$$

$$\leq 2\sqrt{\beta_{t-1}(\delta)} \|x_t\|_{G_{t-1}^{-1}} \tag{19}$$

where $\widetilde{\theta}_t$ is the optimistic parameter estimate, *i.e.* the $\theta \in C_t$ that maximizes the upper confidence bound on the reward of action $x_t$. The first inequality is since $\left(X_t, \widetilde{\theta}_t\right)$ is optimistic, and the last step holds by Cauchy-Schwarz. Using inequality (19) we get that, with probability at least $1 - \delta$, for all $n \geq 0$

$$R_n \leq \sqrt{n \sum_{t=1}^{n} r_t^2} \leq \sqrt{8\beta_n(\delta) n \sum_{t=1}^{n} \|x_t\|_{G_{t-1}^{-1}}}$$

$$\leq 4\sqrt{\frac{ndX^2 \log(1 + nX^2/\lambda d)}{\lambda \log(1 + X^2/\lambda)}} \left(\lambda^{1/2} S + \sigma \sqrt{2 \log(1/\delta) + d \log(1 + nX^2/(\lambda d))}\right)$$

where the last step follow from Lemma C.1. $\qquad\qquad\square$

## F.2 Non-stationary linear bandits

In this section, we study linear stochastic bandits in the non-stationary setting. We provide an experimental study of this setup in Section G. We now turn to the setting of *non-stationary stochastic linear bandits*, where the target parameter is varying with time: $\theta_* = \theta_*(t) \in \mathbb{R}^d$, assuming that $\sum_{s=1}^{T-1} \|\theta_*(s) - \theta_*(s+1)\|_2 \leq B_T$.

One of the optimal algorithms in this setting is `D-LinUCB` of [19], it defines $\theta_t$ as

$$\theta_t = \underset{\theta \in \mathbb{R}^d}{\arg\min} \sum_{s=1}^{t} \gamma^{t-s} (y_s - \langle x_s, \theta \rangle)^2 + \lambda/2 \|\theta\|_2^2.$$

`D-LinUCB` proceeds as follows:

---
**Algorithm 4:** `D-LinUCB`

---
**Input:** $\delta, \sigma, \lambda, X, S, \gamma > 0$, dimension $d \in \mathbb{N}^*$.
**Initialization:** $b = 0_{\mathbb{R}^d}$, $V = \lambda I_d$, $\widetilde{V} = \lambda I_d$, $\theta = 0_{\mathbb{R}^d}$
**for** $t \geq 1$ **do**

> Receive $\mathcal{X}$, compute $\beta_{t-1} = \sqrt{\lambda} S + \sigma \sqrt{2 \log\left(\frac{1}{\delta}\right) + d \log\left(1 + \frac{X^2(1 - \gamma^{2(t-1)})}{\lambda d(1 - \gamma^2)}\right)}$
>
> **for** $a \in \mathcal{X}$ **do**
>> Compute $\text{UCB}(a) = a^\top \theta + \beta_{t-1} \sqrt{a^\top V^{-1} \widetilde{V} V^{-1} a}$
>
> $A_t = \arg\max_a (\text{UCB}(a))$
> **Play action** $A_t$ and **receive reward** $X_t$
> **Updating phase**: $V = \gamma V + x_t x_t^\top + (1 - \gamma)\lambda I_d$, $\widetilde{V} = \gamma^2 \widetilde{V} + x_t x_t^\top + (1 - \gamma^2)\lambda I_d$
> $\qquad\qquad b = \gamma b + Y_t X_t$, $\theta = V^{-1} b$

---

We recall the regret bound of standard `D-LinUCB` .

**Theorem F.1.** *(Theorem 3 of Russac et al. [19]) Assuming that $\sum_{s=1}^{T-1} \|\theta_*(s) - \theta_*(s+1)\|_2 \leq B_T$ and $\forall x \in \mathcal{X}, t \geq 1 : \langle x, \theta_t \rangle \leq 1$, the regret of the D-LinUCB algorithm is bounded for all $\gamma, \delta \in (0, 1)$ and integer $D \geq 1$, with probability at least $1 - \delta$, by:*

$$R_T^r \leq 2XDB_T + \frac{4X^3 S}{\lambda} \frac{\gamma^D}{1 - \gamma} T + 2\sqrt{2}\beta_T \sqrt{dT} \times \sqrt{T \log(1/\gamma) + \log\left(1 + \frac{X^2}{d\lambda(1 - \gamma)}\right)},$$

*where $\beta_T$ is the width of the confidence interval for $\theta_*(T)$.*

Now we introduce `D-LinUCB` $^f$, which uses the forward algorithm and defines an action dependent $\theta_t$ as:

$$\underset{\theta \in \mathbb{R}^d}{\arg\min} \sum_{s=1}^{t} \gamma^{t-s} (y_s - \langle x_s, \theta \rangle)^2 + \lambda/2 \|\theta\|_2^2 + \langle x, \theta \rangle^2. \tag{20}$$

**Theorem F.2.** *Assuming that $\sum_{s=1}^{T-1} \|\theta_*(s) - \theta_*(s+1)\|_2 \leq B_T$, the regret of the `D-LinUCB`$^f$ is bounded for all $\gamma, \delta \in (0, 1)$ and integer $D \geq 1$, with probability at least $1 - \delta$, by*

$$R_T^f \leq 2XDB_T + \frac{4X^3S}{\lambda}\frac{\gamma^D}{1-\gamma}T + 2\beta_T\sqrt{dT}\sqrt{T\log(1/\gamma) + \log\left(1 + \frac{(2-\gamma)X^2}{d\lambda(1-\gamma)}\right)}.$$

*Proof.* This result is again a modification of the original proof consisting in bounding the sum of the actions' norms differently. Let us recall the notations $V_t = \sum_{s=1}^{t} w_s x_s x_s^\top + \lambda_t I_d + xx^\top$ and $\tilde{V}_t = \sum_{s=1}^{t} w_s^2 x_s x_s^\top + \mu_t I_d + xx^\top$. To summarize the difference of this analysis -that no longer requires a bounded rewards assumption- at the step where we bound the sum of actions' norms, we replace Proposition 4 of Russac et al. [19]:

$$\sum_{t=1}^{T} \min\left(1, \|x_t\|_{V_{t-1}^{-1}\tilde{V}_{t-1}V_{t-1}^{-1}}^2\right) \leq 2\sum_{t=1}^{T} \log\left(1 + \gamma^{-t}\|x_t\|_{V_{t-1}^{-1}}^2\right) \leq 2\log\left(\frac{\det(V_T)}{\lambda^d}\right),$$

that requires the predictions to lie in the same range as the rewards with this inequality for `D-LinUCB`$^f$

$$\sum_{t=1}^{T} \|x_t\|_{V_t^{-1}\tilde{V}_tV_t^{-1}}^2 \leq \sum_{t=1}^{T} \log\left(1 + \gamma^{-t}\|x_t\|_{V_t^{-1}}^2\right) \leq \log\left(\frac{\det(V_T)}{\lambda^d}\right).$$

We don't provide the full proof of this result as it is cumbersome and not of special interest for our purposes since it is similar to the analysis for `D-LinUCB` except for the inequality above. $\square$

**Remark 6.** *This result is fascinating as it first allows to remove an unnecessary assumption, and further yields a better bound than `D-LinUCB`$^r$ which suffers the factor $\frac{X\sqrt{2}}{\lambda\log(1+X/\lambda)}$ in its last regret term without assumption 1.*

## G   Experiments

**Experimental details and instructions:**   The experiments were run on a personal laptop with Intel Core i7-8665U, CPU 1.90GHz × 8. Code for the experiments for online regression and linear bandits is provided in the files "OnlineRegression.ipynb" and "LinearBanditsCode.ipynb". For the experiments of non-stationary linear bandits that we present next, we used an existing code from the Github page of Russac et al. [19] and we added an implementation of `D-LinUCB`$^f$ to compare with previous algorithms, this can be seen in the "WeightedLinearBandits" folder in which "D-LinUCB Forward_class.py" is our new algorithm; experiments for this setting can be run from the two ipynb files in the Experiments sub-folder.

**Experiments for non-stationary linear bandits:**   We now reproduce the experiments of [19] for non-stationary linear bandits, and add `D-LinUCB`$^f$ to the pool of algorithms. We first simulate an *abruptly* changing environment of dimension 2 with 3 changes: for $t < 10^3 : \theta_* = (1, 0)$; for $10^3 \leq t \leq 2.10^3 : \theta_* = (-1, 0)$; for $2.10^3 < t < 3.10^3 : \theta_* = (0, 1)$; for $t > 3.10^3 : \theta_* = (0, -1)$. We observe in Fig. 4a that both variants of `D-LinUCB` compare on par. Here `LinUCB-OR` denotes an oracle knowing the change points.

Second, we simulate a slowly changing environment where the parameter $\theta_*$ starts at $(1, 0)$ and moves counter-clockwise on the unit-circle up to the position $(0, 1)$ in $3.10^3$ steps then remains there, $B_T = 1.57$. We see the results in Fig. 4b, where we notice that in this setting as well, `D-LinUCB`$^f$ has very similar performance to standard `D-LinUCB`.

**Remark 7.** *In both experiments, we also reported the performances of `SW-LinUCB`, that is alternative version to `D-LinUCB`. `SW-LinUCB` is better suited for abrupt changes while `D-LinUCB` is better suited for slow changes.*

Note that we added these final experiments to demonstrate the competitiveness of algorithms that use forward regression against their ridge counterparts in the same settings that were used by previous works. While we could have specified specific parameters to illustrate the robustness to regularization of algorithms that incorporate the forward algorithm; we estimate that the experiments presented in the main text already fulfilled this objective. Again, the purpose here is to show that using the forward algorithm improves the theoretical guarantees without deteriorating the performance.

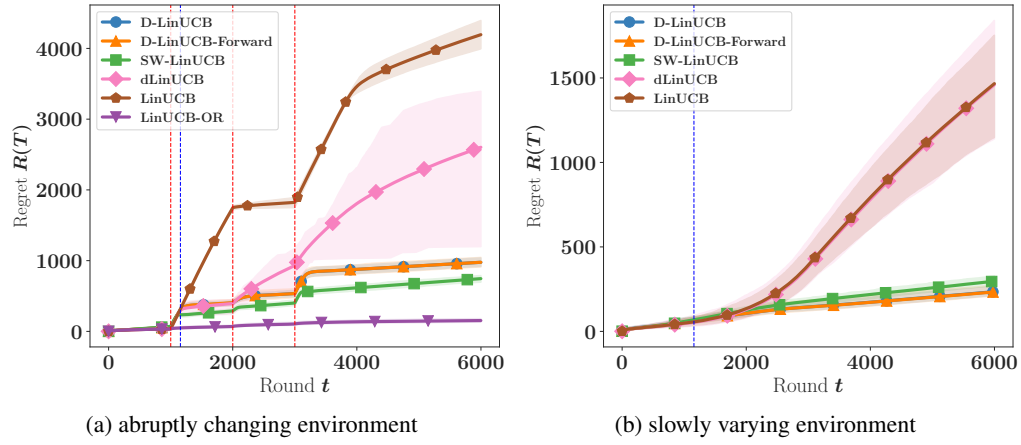(a) abruptly changing environment        (b) slowly varying environment

Figure 4: Performance of several algorithms in an non-stationary environments, averaged over 100 runs, shaded areas represent one standard deviation.