



Second-order kernel online convex optimization with adaptive sketching

Daniele Calandriello, Alessandro Lazaric, Michal Valko
Sequel, Inria Lille – Nord Europe

ICML 2017

ICML, August 2017

Linear Online Convex Optimization (LOCO)

Online game between learner and adversary, at each round $t \in [T]$

- 1 the **adversary** reveals a new point $\mathbf{x}_t \in \mathcal{X}$
- 2 the learner chooses a function $f_{\mathbf{w}_t}$ and predicts $f_{\mathbf{w}_t}(\mathbf{x}_t) = \mathbf{x}_t^\top \mathbf{w}_t$,
- 3 the adversary reveals the **curved** loss ℓ_t ,
- 4 the learner suffers $\ell_t(\mathbf{x}_t^\top \mathbf{w}_t)$ and observes the associated gradient \mathbf{g}_t .

Kernel Online Convex Optimization (KOCO)

Online game between learner and adversary, at each round $t \in [T]$

- 1 the **adversary** reveals a new point $\mathbf{x}_t \in \mathcal{X}$
- 2 the learner chooses a function $f_{\mathbf{w}_t}$ and predicts $f_{\mathbf{w}_t}(\mathbf{x}_t) = \mathbf{x}_t^\top \mathbf{w}_t$,
- 3 the adversary reveals the **curved** loss ℓ_t ,
- 4 the learner suffers $\ell_t(\mathbf{x}_t^\top \mathbf{w}_t)$ and observes the associated gradient \mathbf{g}_t .

Kernel

$\varphi(\cdot) : \mathcal{X} \rightarrow \mathcal{H}$ is the **high-dimensional** (possibly infinite) map

Kernel Online Convex Optimization (KOCO)

Online game between learner and adversary, at each round $t \in [T]$

- 1 the **adversary** reveals a new point $\varphi(\mathbf{x}_t) = \phi_t \in \mathcal{H}$
- 2 the learner chooses a function $f_{\mathbf{w}_t}$ and predicts $f_{\mathbf{w}_t}(\mathbf{x}_t) = \varphi(\mathbf{x}_t)^\top \mathbf{w}_t$,
- 3 the adversary reveals the **curved** loss ℓ_t ,
- 4 the learner suffers $\ell_t(\phi_t^\top \mathbf{w}_t)$ and observes the associated gradient \mathbf{g}_t .

Kernel

$\varphi(\cdot) : \mathcal{X} \rightarrow \mathcal{H}$ is the **high-dimensional** (possibly infinite) map

Kernel Online Convex Optimization (KOCO)

Online game between learner and adversary, at each round $t \in [T]$

- 1 the **adversary** reveals a new point $\varphi(\mathbf{x}_t) = \phi_t \in \mathcal{H}$
- 2 the learner chooses a function $f_{\mathbf{w}_t}$ and predicts $f_{\mathbf{w}_t}(\mathbf{x}_t) = \varphi(\mathbf{x}_t)^\top \mathbf{w}_t$,
- 3 the adversary reveals the **curved** loss ℓ_t ,
- 4 the learner suffers $\ell_t(\phi_t^\top \mathbf{w}_t)$ and observes the associated gradient \mathbf{g}_t .

Kernel

$\varphi(\cdot) : \mathcal{X} \rightarrow \mathcal{H}$ is the **high-dimensional** (possibly infinite) map

$\Phi_t = [\phi_1, \dots, \phi_t]$, $\Phi_t^\top \Phi_t = \mathbf{K}_t$ (kernel trick)

$\mathbf{g}_t = \ell'_t(\phi_t^\top \mathbf{w}_t) \phi_t := \dot{\mathbf{g}}_t \phi_t$

Kernel Online Convex Optimization (KOCO)

Online game between learner and adversary, at each round $t \in [T]$

- 1 the **adversary** reveals a new point $\varphi(\mathbf{x}_t) = \phi_t \in \mathcal{H}$
- 2 the learner chooses a function $f_{\mathbf{w}_t}$ and predicts $f_{\mathbf{w}_t}(\mathbf{x}_t) = \varphi(\mathbf{x}_t)^\top \mathbf{w}_t$,
- 3 the adversary reveals the **curved** loss ℓ_t ,
- 4 the learner suffers $\ell_t(\phi_t^\top \mathbf{w}_t)$ and observes the associated gradient \mathbf{g}_t .

Kernel

$\varphi(\cdot) : \mathcal{X} \rightarrow \mathcal{H}$ is the **high-dimensional** (possibly infinite) map

$\Phi_t = [\phi_1, \dots, \phi_t]$, $\Phi_t^\top \Phi_t = \mathbf{K}_t$ (kernel trick)

$\mathbf{g}_t = \ell'_t(\phi_t^\top \mathbf{w}_t) \phi_t := \dot{\mathbf{g}}_t \phi_t$

Optimization to minimize **regret** $R(\mathbf{w}) = \sum_{t=1}^T \ell_t(\phi_t^\top \mathbf{w}_t) - \ell_t(\phi_t^\top \mathbf{w})$
and **compete** with **best-in-hindsight** $\mathbf{w}^* := \arg \min_{\mathbf{w} \in \mathcal{H}} \sum_{t=1}^T \ell_t(\phi_t^\top \mathbf{w})$

Kernel Online Convex Optimization (KOCO)



convex

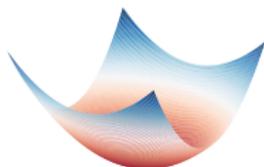
First order (GD) Zinkevich 2003, Kivinen et al. 2004:

- ↳ $\mathcal{O}(d)/\mathcal{O}(t)$ time/space per-step but **slow rate**
approximation avoids $\mathcal{O}(t)$ runtime dependency
 - ↳ but introduce approximation error (potentially $\mathcal{O}(T)$ regret)

Kernel Online Convex Optimization (KOCO)



convex



strongly convex

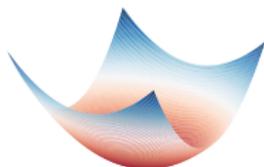
First order (GD) Hazan, Rakhlin, et al. 2008:

↳ $\mathcal{O}(d)/\mathcal{O}(t)$ time/space per-step and fast rate
but not satisfied in practice (e.g. $(y_t - \phi_t^\top \mathbf{w}_t)^2$)

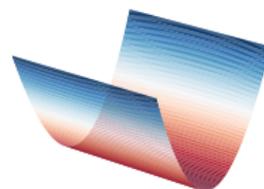
Kernel Online Convex Optimization (KOCO)



convex



strongly convex



σ curved

First order (GD)

↳ $\mathcal{O}(d)/\mathcal{O}(t)$ time/space per-step but slow rate

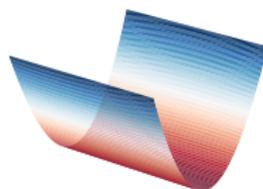
Kernel Online Convex Optimization (KOCO)



convex



strongly convex



σ curved

First order (GD)

↳ $\mathcal{O}(d)/\mathcal{O}(t)$ time/space per-step but slow rate

Second order Hazan, Kalai, et al. 2006, Zhdanov and Kalnishkan 2010:

↳ fast rate but $\mathcal{O}(d^2)/\mathcal{O}(t^2)$ time/space per-step

fast approximations for linear case Luo et al. 2016

no approximate methods for kernel case

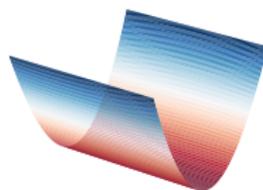
Kernel Online Convex Optimization (KOCO)



convex



strongly convex



σ curved

First order (GD)

↳ $\mathcal{O}(d)/\mathcal{O}(t)$ time/space per-step but slow rate

Second order Hazan, Kalai, et al. 2006, Zhdanov and Kalnishkan 2010:

↳ fast rate but $\mathcal{O}(d^2)/\mathcal{O}(t^2)$ time/space per-step

fast approximations for linear case Luo et al. 2016

no approximate methods for kernel case

How to reduce computational cost without losing fast rate?

Second-Order KOCO (Kernel-Online Newton Step)

Second-Order Gradient Descent

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \mathbf{A}_t^{-1} \mathbf{g}_t, \quad \mathbf{A}_t = \mathbf{A}_{t-1} + \sigma \mathbf{g}_t \mathbf{g}_t^\top, \quad \mathbf{A}_0 = \alpha \mathbf{I}$$

$$\mathbf{A}_t^{-1} = \left(\begin{array}{c} \text{Vertical bars} \\ \times \\ \text{Horizontal bars} \\ + \\ \text{Diagonal matrix} \end{array} \right)^{-1}$$

$$= \begin{array}{c} \text{Diagonal matrix} \\ + \\ \text{Vertical bars} \left(\begin{array}{c} \text{Grid} \\ + \\ \text{Diagonal matrix} \end{array} \right)^{-1} \text{Horizontal bars} \end{array}$$

The diagram illustrates the matrix inversion process for \mathbf{A}_t^{-1} . The top equation shows the matrix as the inverse of the sum of a product of a vertical matrix and a horizontal matrix, and a diagonal matrix. The bottom equation shows the matrix as the sum of a diagonal matrix and a product of a vertical matrix, a grid matrix, a diagonal matrix, and a horizontal matrix.

Second-Order KOCO (Kernel-Online Newton Step)

Second-Order Gradient Descent

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \mathbf{A}_t^{-1} \mathbf{g}_t, \quad \mathbf{A}_t = \mathbf{A}_{t-1} + \sigma \mathbf{g}_t \mathbf{g}_t^\top, \quad \mathbf{A}_0 = \alpha \mathbf{I}$$

$$\begin{aligned} R(\mathbf{w}) &\leq \mathcal{O} \left(\sum_{t=1}^T \mathbf{g}_t^\top \mathbf{A}_t^{-1} \mathbf{g}_t \right) \leq \mathcal{O} \left(L \sum_{t=1}^T \phi_t^\top (\Phi_t \Phi_t^\top + \alpha \mathbf{I})^{-1} \phi_t \right) \\ &\leq \mathcal{O}(\log(\text{Det}(\mathbf{K}_T + \alpha \mathbf{I}))) \end{aligned}$$

Second-Order KOCO (Kernel-Online Newton Step)

Second-Order Gradient Descent

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \mathbf{A}_t^{-1} \mathbf{g}_t, \quad \mathbf{A}_t = \mathbf{A}_{t-1} + \sigma \mathbf{g}_t \mathbf{g}_t^\top, \quad \mathbf{A}_0 = \alpha \mathbf{I}$$

$$\begin{aligned} R(\mathbf{w}) &\leq \mathcal{O} \left(\sum_{t=1}^T \mathbf{g}_t^\top \mathbf{A}_t^{-1} \mathbf{g}_t \right) \leq \mathcal{O} \left(L \sum_{t=1}^T \phi_t^\top (\Phi_t \Phi_t^\top + \alpha \mathbf{I})^{-1} \phi_t \right) \\ &\leq \mathcal{O}(\log(\text{Det}(\mathbf{K}_T + \alpha \mathbf{I}))) \leq ? \end{aligned}$$

Second-Order KOCO (Kernel-Online Newton Step)

Second-Order Gradient Descent

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \mathbf{A}_t^{-1} \mathbf{g}_t, \quad \mathbf{A}_t = \mathbf{A}_{t-1} + \sigma \mathbf{g}_t \mathbf{g}_t^\top, \quad \mathbf{A}_0 = \alpha \mathbf{I}$$

$$\begin{aligned} R(\mathbf{w}) &\leq \mathcal{O} \left(\sum_{t=1}^T \mathbf{g}_t^\top \mathbf{A}_t^{-1} \mathbf{g}_t \right) \leq \mathcal{O} \left(L \sum_{t=1}^T \underbrace{\phi_t^\top (\boldsymbol{\Phi}_t \boldsymbol{\Phi}_t^\top + \alpha \mathbf{I})^{-1} \phi_t}_{\text{Ridge Leverage Score } \tau_{t,t}} \right) \\ &\leq \mathcal{O}(\log(\text{Det}(\mathbf{K}_T + \alpha \mathbf{I}))) \leq ? \end{aligned}$$

Second-Order KOCO (Kernel-Online Newton Step)

Second-Order Gradient Descent

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \mathbf{A}_t^{-1} \mathbf{g}_t, \quad \mathbf{A}_t = \mathbf{A}_{t-1} + \sigma \mathbf{g}_t \mathbf{g}_t^\top, \quad \mathbf{A}_0 = \alpha \mathbf{I}$$

$$\begin{aligned} R(\mathbf{w}) &\leq \mathcal{O} \left(\sum_{t=1}^T \mathbf{g}_t^\top \mathbf{A}_t^{-1} \mathbf{g}_t \right) \leq \mathcal{O} \left(L \sum_{t=1}^T \underbrace{\phi_t^\top (\boldsymbol{\Phi}_t \boldsymbol{\Phi}_t^\top + \alpha \mathbf{I})^{-1} \phi_t}_{\text{Ridge Leverage Score } \tau_{t,t}} \right) \\ &\leq \mathcal{O}(\log(\text{Det}(\mathbf{K}_T + \alpha \mathbf{I}))) \leq ? \end{aligned}$$

Lemma: $\mathcal{O}(\log(\text{Det}(\mathbf{K}_T + \alpha \mathbf{I}))) \leq 2\mathbf{d}_{\text{eff}}^T(\alpha) \log(\mathbf{T}/\alpha)$.

Effective Dimension

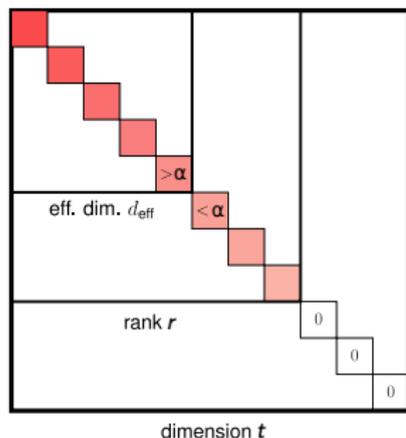
Formally $d_{\text{eff}}^T(\alpha)$ is an α soft-thresholded version of the rank defined as

$$d_{\text{eff}}^T(\alpha) = \text{Tr} \left(\mathbf{K}_T (\mathbf{K}_T + \alpha \mathbf{I})^{-1} \right) = \sum_{t=1}^T \frac{\lambda_t(\mathbf{K}_T)}{\lambda_t(\mathbf{K}_T) + \alpha} \leq \text{Rank}(\mathbf{K}_T) = r$$

Effective Dimension

Formally $d_{\text{eff}}^T(\alpha)$ is an α soft-thresholded version of the rank defined as

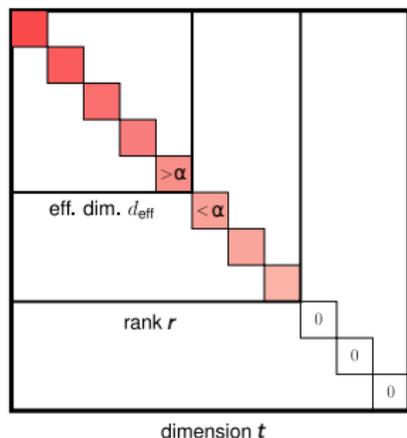
$$d_{\text{eff}}^T(\alpha) = \text{Tr} \left(\mathbf{K}_T (\mathbf{K}_T + \alpha \mathbf{I})^{-1} \right) = \sum_{t=1}^T \frac{\lambda_t(\mathbf{K}_T)}{\lambda_t(\mathbf{K}_T) + \alpha} \leq \text{Rank}(\mathbf{K}_T) = r$$



Effective Dimension

Formally $d_{\text{eff}}^T(\alpha)$ is an α soft-thresholded version of the rank defined as

$$d_{\text{eff}}^T(\alpha) = \text{Tr} \left(\mathbf{K}_T (\mathbf{K}_T + \alpha \mathbf{I})^{-1} \right) = \sum_{t=1}^T \frac{\lambda_t(\mathbf{K}_T)}{\lambda_t(\mathbf{K}_T) + \alpha} \leq \text{Rank}(\mathbf{K}_T) = r$$

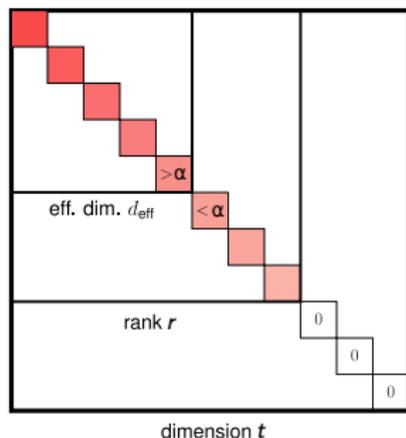


Intuitively, it quantifies the number of **relevant orthogonal directions** played by the adversary.

Effective Dimension

Formally $d_{\text{eff}}^T(\alpha)$ is an α soft-thresholded version of the rank defined as

$$d_{\text{eff}}^T(\alpha) = \text{Tr} \left(\mathbf{K}_T (\mathbf{K}_T + \alpha \mathbf{I})^{-1} \right) = \sum_{t=1}^T \frac{\lambda_t(\mathbf{K}_T)}{\lambda_t(\mathbf{K}_T) + \alpha} \leq \text{Rank}(\mathbf{K}_T) = r$$

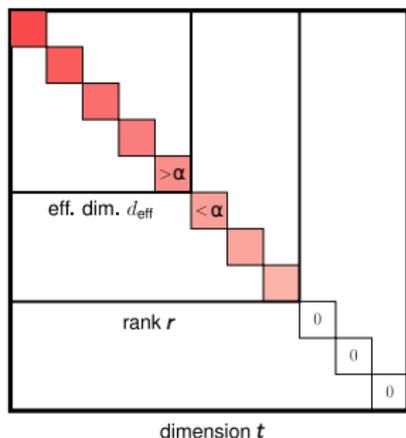


A direction (eigenvector) is relevant if its importance (eigenvalue) is larger than the regularization α

Effective Dimension

Formally $d_{\text{eff}}^T(\alpha)$ is an α soft-thresholded version of the rank defined as

$$d_{\text{eff}}^T(\alpha) = \text{Tr} \left(\mathbf{K}_T (\mathbf{K}_T + \alpha \mathbf{I})^{-1} \right) = \sum_{t=1}^T \frac{\lambda_t(\mathbf{K}_T)}{\lambda_t(\mathbf{K}_T) + \alpha} \leq \text{Rank}(\mathbf{K}_T) = r$$



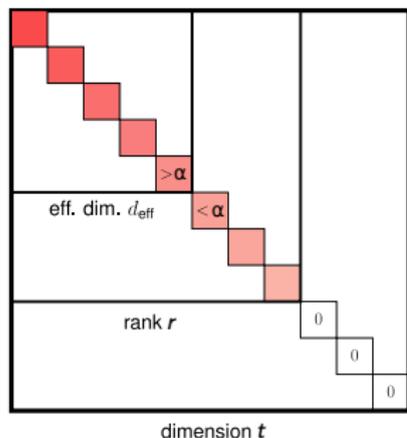
If all ϕ_t are orthogonal

$$d_{\text{eff}}^T(\alpha) \sim T/\alpha$$

Effective Dimension

Formally $d_{\text{eff}}^T(\alpha)$ is an α soft-thresholded version of the rank defined as

$$d_{\text{eff}}^T(\alpha) = \text{Tr} \left(\mathbf{K}_T (\mathbf{K}_T + \alpha \mathbf{I})^{-1} \right) = \sum_{t=1}^T \frac{\lambda_t(\mathbf{K}_T)}{\lambda_t(\mathbf{K}_T) + \alpha} \leq \text{Rank}(\mathbf{K}_T) = r$$



If all ϕ_t are orthogonal

$$d_{\text{eff}}^T(\alpha) \sim T/\alpha$$

If all ϕ_t come from a bounded distribution or a finite set and $\alpha = 1$ then

$$d_{\text{eff}}^T(1) \sim \mathcal{O}(1) \leq r$$

is constant in T

Goal

How to maintain $d_{\text{eff}}^T(\alpha) \log(T)$ regret and reduce computational costs?

Goal

How to maintain $d_{\text{eff}}^T(\alpha) \log(T)$ regret and reduce computational costs?

Computation scales with number of vectors $\mathbf{g}_t \mathbf{g}_t$ added to \mathbf{A}_t

↳ skip some of the additions

Goal

How to maintain $d_{\text{eff}}^T(\alpha) \log(T)$ regret and reduce computational costs?

Computation scales with number of vectors $\mathbf{g}_t \mathbf{g}_t$ added to \mathbf{A}_t

↳ skip some of the additions

Regret is large when $\tau_{t,t}$ is large

↳ Update \mathbf{A}_t only when $\tau_{t,t}$ is large

Goal

How to maintain $d_{\text{eff}}^T(\alpha) \log(T)$ regret and reduce computational costs?

Computation scales with number of vectors $\mathbf{g}_t \mathbf{g}_t$ added to \mathbf{A}_t

↳ skip some of the additions

Regret is large when $\tau_{t,t}$ is large

↳ Update \mathbf{A}_t only when $\tau_{t,t}$ is large

Computing $\tau_{t,t}$ is expensive

↳ Propose KERNEL ONLINE ROW SAMPLING (KORS)
to approximate $\tau_{t,t}$ efficiently

Naive Approach

Unbiased estimator:

$$\tilde{\mathbf{A}}_t = \tilde{\mathbf{A}}_{t-1} + (\mathbb{I}\{\text{coin flip w.p. } p_t\} / p_t) \sigma \mathbf{g}_t \mathbf{g}_t^\top$$

with $p_t \propto \tilde{\tau}_{t,t}$

Naive Approach

Unbiased estimator:

$$\tilde{\mathbf{A}}_t = \tilde{\mathbf{A}}_{t-1} + (\mathbb{I}\{\text{coin flip w.p. } p_t\} / \mathbf{p}_t) \sigma \mathbf{g}_t \mathbf{g}_t^\top$$

with $p_t \propto \tilde{\tau}_{t,t}$

$$\tilde{\mathbf{A}}_t^{-1} = \left(\begin{array}{|c|} \hline \text{red} \\ \hline \text{white} \\ \hline \text{orange} \\ \hline \end{array} \times \begin{array}{|c|} \hline \text{red} \\ \hline \text{white} \\ \hline \text{orange} \\ \hline \end{array} + \begin{array}{|c|} \hline \text{diag} \\ \hline \end{array} \right)^{-1}$$

$$= \begin{array}{|c|} \hline \text{diag} \\ \hline \end{array} + \begin{array}{|c|} \hline \text{red} \\ \hline \text{orange} \\ \hline \end{array} \left(\begin{array}{|c|} \hline \text{red} \\ \hline \text{orange} \\ \hline \end{array} + \begin{array}{|c|} \hline \text{diag} \\ \hline \end{array} \right)^{-1} \begin{array}{|c|} \hline \text{red} \\ \hline \text{orange} \\ \hline \end{array}$$

d_{eff} d_{eff} ∞

Naive Approach

Unbiased estimator:

$$\tilde{\mathbf{A}}_t = \tilde{\mathbf{A}}_{t-1} + (\mathbb{I}\{\text{coin flip w.p. } p_t\} / p_t) \sigma \mathbf{g}_t \mathbf{g}_t^\top$$

with $p_t \propto \tilde{\tau}_{t,t}$

Pros:

w.h.p. $\tilde{\mathbf{A}}_t$ updated only $d_{\text{eff}}^T(\alpha) \log^2(T)$ times

$\tilde{O}(d_{\text{eff}}^T(\alpha)^2 + t)$ per-step space/time complexity

Cons:

Expected regret $d_{\text{eff}}^T(\alpha) \log(T)$

The weights $1/p_t \sim 1/\tilde{\tau}_{t,t}$ can be large

↳ large updates $\tilde{\mathbf{A}}_t - \tilde{\mathbf{A}}_{t-1}$ and large variance

Naive Approach

biased estimator:

$$\tilde{\mathbf{A}}_t = \tilde{\mathbf{A}}_{t-1} + (\mathbb{I}\{\text{coin flip w.p. } p_t\}) \sigma \mathbf{g}_t \mathbf{g}_t^\top$$

with $p_t \propto \tilde{\tau}_{t,t}$

Sketched-KONS

SKETCHED-KONS:

$$\tilde{\mathbf{A}}_t = \tilde{\mathbf{A}}_{t-1} + (\mathbb{I}\{\text{coin flip w.p. } p_t\}) \sigma \mathbf{g}_t \mathbf{g}_t^\top$$

with $p_t \propto \max\{\gamma, \tilde{\tau}_{t,t}\}$

$$\tilde{\mathbf{A}}_t^{-1} = \left(\begin{array}{c} \text{vertical bars (red, green, orange)} \\ \times \\ \text{horizontal bars (red, green, orange)} \end{array} + \begin{array}{c} \text{diagonal matrix} \end{array} \right)^{-1}$$

$$= \begin{array}{c} \text{diagonal matrix} \\ + \\ \begin{array}{c} \text{vertical bars (red, green, orange)} \\ \text{diagonal matrix} \\ + \\ \text{diagonal matrix} \end{array}^{-1} \end{array} \begin{array}{c} \text{horizontal bars (red, green, orange)} \end{array}$$

Sketched-KONS

SKETCHED-KONS:

$$\tilde{\mathbf{A}}_t = \tilde{\mathbf{A}}_{t-1} + (\mathbb{I}\{\text{coin flip w.p. } p_t\}) \sigma \mathbf{g}_t \mathbf{g}_t^\top$$

with $p_t \propto \max\{\gamma, \tilde{\tau}_{t,t}\}$

Theorem: with high probability SKETCHED-KONS

(1) achieves $d_{\text{eff}}^T(\alpha) \log(T) / \max\{\gamma, \tau_{\min}\}$ regret

(2) requires only $\tilde{\mathcal{O}}(d_{\text{eff}}^T(\alpha)^2 + \mathbf{t}^2 \gamma^2 + t)$ per-step space/time

Trade-off $1/\gamma$ increase in regret for γ^2 space/time improvement

Still dependent on t

What next?

Can we get rid of dependency on t without losing fast rates?

Skipping updates not enough:

↳ we propose a counterexample

Different approaches?

Learn how to remove old $\mathbf{g}_s \mathbf{g}_s^\top$ from \mathbf{A}_t ?

Instead of approximating \mathbf{A}_t , approximate ϕ_t

↳ Random feature not strong enough yet
(Avron et al. ICML 17 for batch setting)
Ongoing work using RLS sampling and Nyström embeddings

Bibliography I

-  Elad Hazan, Adam Kalai, Satyen Kale, and Amit Agarwal. “Logarithmic regret algorithms for online convex optimization”. In: [Conference on Learning Theory](#). Springer, 2006, pp. 499–513.
-  Elad Hazan, Alexander Rakhlin, and Peter L Bartlett. “Adaptive online gradient descent”. In: [Advances in Neural Information Processing Systems](#). 2008, pp. 65–72.
-  J. Kivinen, A.J. Smola, and R.C. Williamson. “Online Learning with Kernels”. en. In: [IEEE Transactions on Signal Processing](#) 52.8 (Aug. 2004). (Visited on 02/18/2017).
-  Haipeng Luo, Alekh Agarwal, Nicolo Cesa-Bianchi, and John Langford. “Efficient second-order online learning via sketching”. In: [Neural Information Processing Systems](#). 2016.

Bibliography II



Fedor Zhdanov and Yuri Kalnishkan. “An Identity for Kernel Ridge Regression”. en. In: Algorithmic Learning Theory. Ed. by Springer. Lecture Notes in Computer Science. Oct. 2010, pp. 405–419.



Martin Zinkevich. “Online Convex Programming and Generalized Infinitesimal Gradient Ascent”. In: International Conference on Machine Learning. 2003, pp. 928–936.