



# Efficient second-order online kernel learning with adaptive embedding

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

COLT 2017

COLT impromptu talks, July 2017, La France avance!

## Online Kernel Learning (OKL)

**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  $\ell_t$ ,
- 4 the learner suffers  $\ell_t(\phi_t^\top \mathbf{w}_t)$  and observes the associated gradient  $\mathbf{g}_t$ .

# Online Kernel Learning (OKL)

**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  $\ell_t$ ,
- 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)

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

# Online Kernel Learning (OKL)

**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  $\ell_t$ ,
- 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)

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

**Learning** to minimize **regret**  $R(\mathbf{w}) = \sum_{t=1}^T \ell_t(\phi_t \mathbf{w}_t) - \ell_t(\phi_t \mathbf{w})$

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

## Curved losses?

We assume the losses  $l_t$  are scalar Lipschitz

$$|l'_t(z)| \leq L \text{ whenever } |z| \leq C$$

and curved

$$l_t(\phi_t^\top \mathbf{w}) \geq l_t(\phi_t^\top \mathbf{u}) + \nabla l_t(\phi_t^\top \mathbf{u})^\top (\mathbf{w} - \mathbf{u}) + \sigma (\nabla l_t(\phi_t^\top \mathbf{u})^\top (\mathbf{w} - \mathbf{u}))^2.$$

## Curved losses?

We assume the losses  $l_t$  are scalar Lipschitz

$$|l'_t(z)| \leq L \text{ whenever } |z| \leq C$$

and curved

$$l_t(\phi_t^T \mathbf{w}) \geq l_t(\phi_t^T \mathbf{u}) + \nabla l_t(\phi_t^T \mathbf{u})^T (\mathbf{w} - \mathbf{u}) + \sigma (\nabla l_t(\phi_t^T \mathbf{u})^T (\mathbf{w} - \mathbf{u}))^2.$$

You already use curved losses

weaker than strong convexity

↳ strongly convex only along  $\nabla l_t(\phi_t^T \mathbf{u})$

satisfied by exp-concave losses

↳ squared loss, squared hinge-loss

## Second-Order OKL (Kernel Online Newton Step)

Second-Order Gradient Descent

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

If  $\varphi(\mathbf{x}) = \mathbf{x} : \mathbb{R}^d \rightarrow \mathbb{R}^d$  is the identity, Online Newton Step Hazan et al. 2006

↳  $\mathcal{O}(d^2)$  time/space per-step (Bottleneck: storing and inverting  $\mathbf{X}_t^\top \mathbf{X}_t$ )

$$R(\mathbf{w}^*) \leq \alpha \|\mathbf{w}^* - \mathbf{w}_0\|_2^2 + d \log(\mathbf{T})$$

Sketched-ONS fast approximation, but only for  $\varphi(\mathbf{x})$  identity and low-rank  $\mathbf{X}_t^\top \mathbf{X}_t$  Luo et al. 2016

If  $\varphi(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathcal{H}$  is arbitrary, Kernel-ONS Calandriello et al. 2017

↳  $\mathcal{O}(t^2)$  time/space per-step (Bottleneck: storing and inverting  $\mathbf{K}_t$ )

$$R(\mathbf{w}^*) \leq \underbrace{\alpha \|\mathbf{w}^* - \mathbf{w}_0\|_2^2}_a + \underbrace{\mathbf{d}_{\text{eff}}^T (\alpha / (\mathbf{L}\sigma)) \log(\mathbf{T})}_b$$

(a) start cost, (b) effective dimension (degrees of freedom) of data

## Effective Dimension

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

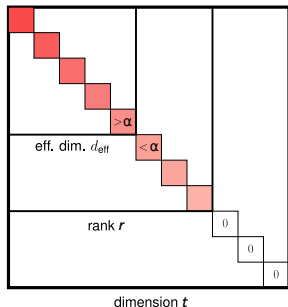
$$d_{\text{eff}}^T\left(\frac{\alpha}{L\sigma}\right) = \text{Tr}\left(\mathbf{K}_T\left(\mathbf{K}_T + \frac{\alpha}{L\sigma}\mathbf{I}\right)^{-1}\right) = \sum_{t=1}^T \frac{\lambda_t}{\lambda_t + \frac{\alpha}{L\sigma}} \leq \text{Rank}(\mathbf{K}_T) = r$$



## Effective Dimension

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

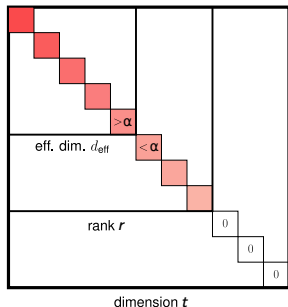
$$d_{\text{eff}}^T\left(\frac{\alpha}{L\sigma}\right) = \text{Tr}\left(\mathbf{K}_T\left(\mathbf{K}_T + \frac{\alpha}{L\sigma}\mathbf{I}\right)^{-1}\right) = \sum_{t=1}^T \frac{\lambda_t}{\lambda_t + \frac{\alpha}{L\sigma}} \leq \text{Rank}(\mathbf{K}_T) = r$$



## Effective Dimension

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

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

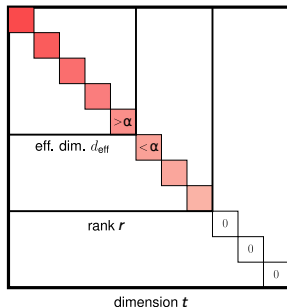


Intuitively, it quantifies the number of relevant orthogonal directions played by the adversary.

## Effective Dimension

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

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

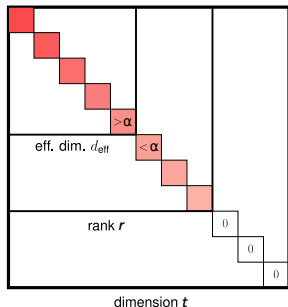


A direction (eigenvector) is relevant if its importance (eigenvalue) is larger than the Lipschitz discounted regularization  $\alpha/(L\sigma)$

## Effective Dimension

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

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



Assume  $\|\phi_t\| = 1$ , then if all  $\phi_t$  are orthogonal and  $\alpha = \sqrt{T}$  then

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

and

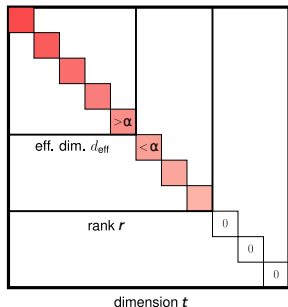
$$R(\mathbf{w}^*) \leq \sqrt{T} + d_{\text{eff}}^T(\sqrt{T}) \log(T) \sim \sqrt{T}$$

recover first order bound.

## Effective Dimension

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

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



If all  $\phi_t$  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  $T$  and

$$R(\mathbf{w}^*) \leq \mathcal{O}(1) + \mathcal{O}(1) \log(T) \sim \log T$$

logarithmic in  $T$ .

## Ideal method

How to achieve adaptive  $d_{\text{eff}}^T(\alpha) \log(T)$  regret  
but without computational complexity depending on  $t$ ?

## Ideal method

How to achieve adaptive  $d_{\text{eff}}^T(\alpha) \log(T)$  regret  
but without computational complexity depending on  $t$ ?

Use approximate second order gradient in  $\mathcal{H}$  Calandriello et al. 2017

↳  $d_{\text{eff}}^T(\alpha) \log(T)$  regret, but runtime still depends on  $t$

## Ideal method

How to achieve adaptive  $d_{\text{eff}}^T(\alpha) \log(T)$  regret  
but without computational complexity depending on  $t$ ?

Use approximate second order gradient in  $\mathcal{H}$  Calandriello et al. 2017

↳  $d_{\text{eff}}^T(\alpha) \log(T)$  regret, but runtime still depends on  $t$

Use exact second order gradient in approximate  $\tilde{\mathcal{H}}$

$$\sum_{t=1}^T \ell_t(\phi_t \tilde{\mathbf{w}}_t) - \ell_t(\phi_t \mathbf{w}^*) = \sum_{t=1}^T \underbrace{\ell_t(\phi_t \tilde{\mathbf{w}}_t) - \ell_t(\phi_t \bar{\mathbf{w}})}_a + \underbrace{\ell_t(\phi_t \bar{\mathbf{w}}) - \ell_t(\phi_t \mathbf{w}^*)}_b$$



## Ideal method

How to achieve adaptive  $d_{\text{eff}}^T(\alpha) \log(T)$  regret  
but without computational complexity depending on  $t$ ?

Use approximate second order gradient in  $\mathcal{H}$  Calandriello et al. 2017

↳  $d_{\text{eff}}^T(\alpha) \log(T)$  regret, but runtime still depends on  $t$

Use exact second order gradient in approximate  $\tilde{\mathcal{H}}$

$$\sum_{t=1}^T \ell_t(\phi_t \tilde{\mathbf{w}}_t) - \ell_t(\phi_t \mathbf{w}^*) = \sum_{t=1}^T \underbrace{\ell_t(\phi_t \tilde{\mathbf{w}}_t) - \ell_t(\phi_t \bar{\mathbf{w}})}_a + \underbrace{\ell_t(\phi_t \bar{\mathbf{w}}) - \ell_t(\phi_t \mathbf{w}^*)}_b$$

(a) error between online and batch in  $\tilde{\mathcal{H}}$ :

↳  $d_{\text{eff}}^T(\alpha) \log(T)$  bound using KONS analysis

## Ideal method

How to achieve adaptive  $d_{\text{eff}}^T(\alpha) \log(T)$  regret  
but without computational complexity depending on  $t$ ?

Use approximate second order gradient in  $\mathcal{H}$  Calandriello et al. 2017

↳  $d_{\text{eff}}^T(\alpha) \log(T)$  regret, but runtime still depends on  $t$

Use exact second order gradient in approximate  $\tilde{\mathcal{H}}$

$$\sum_{t=1}^T \ell_t(\phi_t \tilde{\mathbf{w}}_t) - \ell_t(\phi_t \mathbf{w}^*) = \sum_{t=1}^T \underbrace{\ell_t(\phi_t \tilde{\mathbf{w}}_t) - \ell_t(\phi_t \bar{\mathbf{w}})}_a + \underbrace{\ell_t(\phi_t \bar{\mathbf{w}}) - \ell_t(\phi_t \mathbf{w}^*)}_b$$

(a) error between online and batch in  $\tilde{\mathcal{H}}$ :

↳  $d_{\text{eff}}^T(\alpha) \log(T)$  bound using KONS analysis

(b) error between  $\bar{\mathbf{w}}$  best in  $\tilde{\mathcal{H}}$  and  $\mathbf{w}^*$  best in  $\mathcal{H}$ : bound how?

## Subspace approximation error

$\tilde{\mathcal{H}}$  cannot be fixed

↳ the adversary will find orthogonal points and exploit this

## Subspace approximation error

$\tilde{\mathcal{H}}$  cannot be fixed

↳ the adversary will find orthogonal points and exploit this

Use Nyström approximation instead and adapt it online

↳  $\tilde{\mathcal{H}}_t = \text{Span}(\mathcal{I}_t)$  defined using  $m_t$  inducing points  $\mathcal{I}_t = \{\phi_s\}_{s=1}^{m_t}$

If the adversary plays a "sufficiently orthogonal"  $\phi_t$ , add it to  $\mathcal{I}_{t+1}$

## Subspace approximation error

$\tilde{\mathcal{H}}$  cannot be fixed

↳ the adversary will find orthogonal points and exploit this

Use Nyström approximation instead and adapt it online

↳  $\tilde{\mathcal{H}}_t = \text{Span}(\mathcal{I}_t)$  defined using  $m_t$  inducing points  $\mathcal{I}_t = \{\phi_s\}_{s=1}^{m_t}$

If the adversary plays a "sufficiently orthogonal"  $\phi_t$ , add it to  $\mathcal{I}_{t+1}$

$\tilde{\mathcal{H}}_t$  is finite dimensional: runtime independent of  $t$

↳ Easy to embed (project) points as

$$\tilde{\varphi}(\cdot) = \Sigma^{-1} \mathbf{U}^T \Phi_{\mathcal{I}}^T \varphi(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^{m_t}$$

with  $\mathbf{K}_{\mathcal{I}} = \Phi_{\mathcal{I}}^T \Phi_{\mathcal{I}} = \mathbf{U} \Sigma \Sigma \mathbf{U}^T$

$\mathcal{O}(m_t^2)$  time/space cost to run exact KONS in  $\tilde{\mathcal{H}}_t$

## Subspace approximation error

“sufficiently orthogonal” is measured using  $\gamma$ -ridge leverage scores

$$\mathbb{P}(\text{include } \phi_t \text{ in } \mathcal{I}_t) \sim \phi_t^\top \phi_t - \phi_t^\top (\Phi_{\mathcal{I}_{t-1}} \Phi_{\mathcal{I}_{t-1}}^\top + \gamma \mathbf{I})^{-1} \phi_t$$

Also computable in  $m_t^2$  time.

## Subspace approximation error

“sufficiently orthogonal” is measured using  $\gamma$ -ridge leverage scores

$$\mathbb{P}(\text{include } \phi_t \text{ in } \mathcal{I}_t) \sim \phi_t^\top \phi_t - \phi_t^\top (\Phi_{\mathcal{I}_{t-1}} \Phi_{\mathcal{I}_{t-1}}^\top + \gamma \mathbf{I})^{-1} \phi_t$$

Also computable in  $m_t^2$  time.

Guarantees that Calandriello et al. 2017

$$m_t \leq d_{\text{eff}}^t(\gamma) \log^2(T) \quad (\text{space/time}), \quad \mathbf{K}_t - \tilde{\mathbf{K}}_t \preceq \gamma \mathbf{I} \quad (\text{accuracy})$$

## Online/batch error

Use KONS guarantees to bound  $\sum_{t=1}^T \ell_t(\phi_t \tilde{\mathbf{w}}_t) - \ell_t(\phi_t \bar{\mathbf{w}})$  separately for each  $\tilde{\mathcal{H}}_t$  and associated  $\bar{\mathbf{w}}_j$

$$\sum_{t=1}^T \ell_t(\phi_t \tilde{\mathbf{w}}_t) - \ell_t(\phi_t \bar{\mathbf{w}}) \leq J d_{\text{eff}}^T(\alpha/(L\sigma)) \log(T) + \sum_{j=1}^J \underbrace{\alpha \|\bar{\mathbf{w}}_j - \mathbf{w}_{t_j}\|_2^2}_{\text{start costs}}$$



## Online/batch error

Use KONS guarantees to bound  $\sum_{t=1}^T \ell_t(\phi_t \tilde{\mathbf{w}}_t) - \ell_t(\phi_t \bar{\mathbf{w}})$  separately for each  $\tilde{\mathcal{H}}_t$  and associated  $\bar{\mathbf{w}}_j$

$$\sum_{t=1}^T \ell_t(\phi_t \tilde{\mathbf{w}}_t) - \ell_t(\phi_t \bar{\mathbf{w}}) \leq J d_{\text{eff}}^T(\alpha/(L\sigma)) \log(T) + \sum_{j=1}^J \underbrace{\alpha \|\bar{\mathbf{w}}_j - \mathbf{w}_{t_j}\|_2^2}_{\text{start costs}}$$

Every time we change  $\tilde{\mathcal{H}}$  we pay  $\alpha \|\bar{\mathbf{w}}_j - \mathbf{w}_{t_j}\|_2^2$

↳ the **adversary** can influence  $\mathbf{w}_{t_j}$  and make it large

## Online/batch error

Use KONS guarantees to bound  $\sum_{t=1}^T \ell_t(\phi_t \tilde{\mathbf{w}}_t) - \ell_t(\phi_t \bar{\mathbf{w}})$  separately for each  $\tilde{\mathcal{H}}_t$  and associated  $\bar{\mathbf{w}}_j$

$$\sum_{t=1}^T \ell_t(\phi_t \tilde{\mathbf{w}}_t) - \ell_t(\phi_t \bar{\mathbf{w}}) \leq J d_{\text{eff}}^T(\alpha/(L\sigma)) \log(T) + \sum_{j=1}^J \underbrace{\alpha \|\bar{\mathbf{w}}_j - \mathbf{w}_{t_j}\|_2^2}_{\text{start costs}}$$

Every time we change  $\tilde{\mathcal{H}}$  we pay  $\alpha \|\bar{\mathbf{w}}_j - \mathbf{w}_{t_j}\|_2^2$

↳ the adversary can influence  $\mathbf{w}_{t_j}$  and make it large

Reset  $\tilde{\mathbf{w}}_t$  and  $\tilde{\mathbf{A}}_t$  when  $\tilde{\mathcal{H}}_t$  changes

↳ Wasteful, but not too often. At most  $J \leq d_{\text{eff}}^T(\gamma)$  times.

Learning is preserved through  $\tilde{\mathcal{H}}_t$  that always improves

Adaptive doubling trick

## Final regret guarantees

For any **curved** loss

$$R(\mathbf{w}^*) \leq J(\alpha \|\mathbf{w}^*\|_2^2 + d_{\text{eff}}^T \log(\alpha/(L\sigma))) + \frac{\gamma}{\alpha} T$$

## Final regret guarantees

For any **curved** loss

$$R(\mathbf{w}^*) \leq J(\alpha \|\mathbf{w}^*\|_2^2 + d_{\text{eff}}^T \log(\alpha/(L\sigma))) + \frac{\gamma}{\alpha} T$$

Setting  $\gamma = \alpha/T$  removes second term

↳ computational cost is  $\mathcal{O}(d_{\text{eff}}^T(1/T)^2)$ , still small in many cases

## Final regret guarantees

For any **curved** loss

$$R(\mathbf{w}^*) \leq J(\alpha \|\mathbf{w}^*\|_2^2 + d_{\text{eff}}^T \log(\alpha/(L\sigma))) + \frac{\gamma}{\alpha} T$$

Setting  $\gamma = \alpha/T$  removes second term

↳ computational cost is  $\mathcal{O}(d_{\text{eff}}^T(1/T)^2)$ , still small in many cases

For squared loss only and  $\gamma = \alpha$

$$R(\mathbf{w}^*) \leq J(\alpha \|\mathbf{w}^*\|_2^2 + d_{\text{eff}}^T \log(\alpha/(L\sigma))) + J\left(\sum_{t=1}^T \ell_t(\phi_t \mathbf{w}^*) + \alpha \|\mathbf{w}^*\|_2^2\right)$$

Last term  $J(\sum_{t=1}^T \ell_t(\phi_t \mathbf{w}^*) + \alpha \|\mathbf{w}^*\|_2^2)$  replaces  $\frac{\gamma}{\alpha} T$

↳ **regularized** cumulative loss of  $\mathbf{w}^*$

if  $\mathcal{H}$  is good, very small

# Conclusions

Algorithm	cadata $n = 20k, d = 8$			casp $n = 45k, d = 9$		
	Avg. Squared Loss	#SV	Time	Avg. Squared Loss	#SV	Time
FOGD	$0.04097 \pm 0.00015$	30	—	$0.08021 \pm 0.00031$	30	—
NOGD	$0.03983 \pm 0.00018$	30	—	$0.07844 \pm 0.00008$	30	—
<b>PROS-N-KONS</b>	$0.03095 \pm 0.00110$	20	18.59	<b><math>0.06773 \pm 0.00105</math></b>	21	40.73
CON-KONS	<b><math>0.02850 \pm 0.00174</math></b>	19	18.45	<b><math>0.06832 \pm 0.00315</math></b>	20	40.91
B-KONS	$0.03095 \pm 0.00118$	19	18.65	<b><math>0.06775 \pm 0.00067</math></b>	21	41.13
BATCH	$0.02202 \pm 0.00002$	—	—	$0.06100 \pm 0.00003$	—	—

Algorithm	slice $n = 53k, d = 385$			year $n = 463k, d = 90$		
	Avg. Squared Loss	#SV	Time	Avg. Squared Loss	#SV	Time
FOGD	$0.00726 \pm 0.00019$	30	—	$0.01427 \pm 0.00004$	30	—
NOGD	$0.02636 \pm 0.00460$	30	—	$0.01427 \pm 0.00004$	30	—
DUAL-SGD	—	—	—	$0.01440 \pm 0.00000$	100	—
<b>PROS-N-KONS</b>	did not complete	—	—	$0.01450 \pm 0.00014$	149	884.82
CON-KONS	did not complete	—	—	$0.01444 \pm 0.00017$	147	889.42
B-KONS	<b><math>0.00913 \pm 0.00045</math></b>	100	60	<b><math>0.01302 \pm 0.00006</math></b>	100	505.36
BATCH	$0.00212 \pm 0.00001$	—	—	$0.01147 \pm 0.00001$	—	—

# Bibliography



Daniele Calandriello, Alessandro Lazaric, and Michal Valko. “Second-Order Kernel Online Convex Optimization with Adaptive Sketching”. In: [International Conference on Machine Learning](#). 2017.



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.



Haipeng Luo, Alekh Agarwal, Nicolo Cesa-Bianchi, and John Langford. “Efficient second-order online learning via sketching”. In: [Neural Information Processing Systems](#). 2016.