

Robust Online Convex Optimization in the Presence of Outliers

Tim van Erven



UNIVERSITY
OF AMSTERDAM

COLT 2021



Sarah Sachs



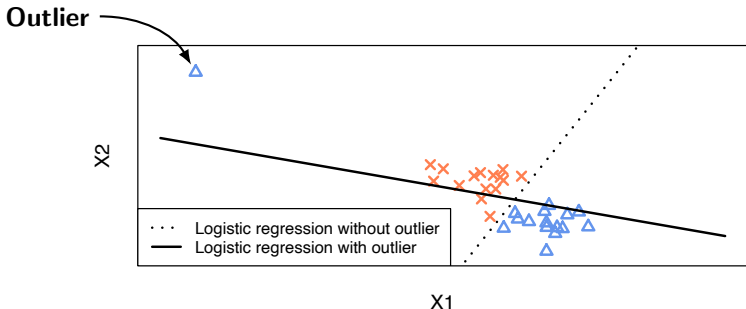
Wouter Koolen



Wojciech Kotłowski

Recruiting: Postdoc position in my group available 2022

Extreme Outliers Can Break Learning



Reasons for outliers:

- ▶ Naturally **heavy-tailed data**
- ▶ A small subset of **malicious users** trying to corrupt data stream
- ▶ Glitches in **cheap sensors**

Heavily studied:

- ▶ In statistics [Tukey, 1959, Huber, 1964], stochastic optimization, etc.
- ▶ But not yet in Online Convex Optimization

Formalizing Robust OCO

Standard OCO setting:

Given convex domain $\mathcal{W} \subset \mathbb{R}^d$ with $\text{diameter}(\mathcal{W}) \leq D$

- 1: **for** $t = 1, 2, \dots, T$ **do**
- 2: Predict w_t in \mathcal{W}
- 3: Observe convex loss function $f_t : \mathcal{W} \rightarrow \mathbb{R}$ with gradient $g_t = \nabla f_t(w_t)$
- 4: **end for**

Robust regret:
$$R_T(\mathbf{u}, \mathcal{S}) = \sum_{t \in \mathcal{S}} (f_t(w_t) - f_t(\mathbf{u}))$$

Challenges:

- ▶ Inliers $\mathcal{S} \subset \{1, \dots, T\}$ **unknown** (chosen by adversary)
- ▶ Bounds **cannot depend on outliers** at all, but must scale with

$$G(\mathcal{S}) = \max_{t \in \mathcal{S}} \|g_t\|.$$

Robustifying Any OCO Algorithm

1. **Any OCO ALG** with regret bound $B_T(G)$ if gradients have length at most G
2. **Top- k Filter**: simple strategy to **filter out large gradients**

Theorem (At most k outliers)

On linear losses, **ALG** + **Top- k Filter** achieves

$$R_T(\mathbf{u}, \mathcal{S}) \leq \underbrace{B_T(2G(\mathcal{S}))}_{\text{Feed ALG gradients}} + 4DG(\mathcal{S})(k+1) \quad \text{for any } \mathcal{S} : T - |\mathcal{S}| \leq k.$$

Feed ALG gradients $\leq 2G(\mathcal{S})$

Robustifying Any OCO Algorithm

1. **Any OCO ALG** with regret bound $B_T(G)$ if gradients have length at most G
2. **Top- k Filter**: simple strategy to **filter out large gradients**

Theorem (At most k outliers)

On linear losses, **ALG** + **Top- k Filter** achieves

$$R_T(\mathbf{u}, \mathcal{S}) \leq B_T(2G(\mathcal{S})) + \underbrace{4DG(\mathcal{S})(k+1)}_{\text{price of robustness}} \quad \text{for any } \mathcal{S} : T - |\mathcal{S}| \leq k.$$

price of robustness = $O(G(\mathcal{S})k)$

Robustifying Any OCO Algorithm

1. **Any OCO ALG** with regret bound $B_T(G)$ if gradients have length at most G
2. **Top- k Filter**: simple strategy to **filter out large gradients**

Theorem (At most k outliers)

On linear losses, **ALG** + **Top- k Filter** achieves

$$R_T(\mathbf{u}, S) \leq B_T(2G(S)) + 4DG(S)(k + 1) \quad \text{for any } S : T - |S| \leq k.$$

Losses	Minimax Robust Regret
General convex	$O(\sqrt{T} + k)$
General convex + i.i.d.	"
Strongly convex	$O(\ln(T) + k)$

Efficient Filtering Approach

Top- k Filter:

- ▶ Maintain list \mathcal{L}_t of $k + 1$ largest gradient lengths seen so far
- ▶ Filter round if $\|g_t\| > 2 \min \mathcal{L}_t$; otherwise pass to ALG

Main Ideas:

1. Never pass ALG gradients $> 2G(\mathcal{S})$:
 - ▶ \mathcal{L}_t contains at least 1 inlier, because at most k outliers
 - ▶ Hence $\min \mathcal{L}_t \leq G(\mathcal{S})$
2. Overhead for filtering is $O(k)$
 - ▶ Every filtered round is also added to \mathcal{L}_t
 - ▶ Therefore $\min \mathcal{L}_t$ (at least) doubles every $k + 1$ filtered rounds
 - ▶ Hence last $k + 1$ filtered rounds dominate

Application: Robustified Online-to-Batch

Huber ϵ -contamination model:

$$P_\epsilon = (1 - \epsilon)P + \epsilon Q$$

Outlier distribution

Distribution of interest

- ▶ $f_t(\mathbf{w}) = f(\mathbf{w}, \xi)$ where $\xi \sim P_\epsilon$
- ▶ Inlier risk: $\text{Risk}_P(\mathbf{w}) = \mathbb{E}_{\xi \sim P}[f(\mathbf{w}, \xi)]$

Application: Robustified Online-to-Batch

Huber ϵ -contamination model:

$$P_\epsilon = (1 - \epsilon)P + \epsilon Q$$

Outlier distribution

Distribution of interest

- ▶ $f_t(\mathbf{w}) = f(\mathbf{w}, \xi)$ where $\xi \sim P_\epsilon$
- ▶ Inlier risk: $\text{Risk}_P(\mathbf{w}) = \mathbb{E}_{\xi \sim P}[f(\mathbf{w}, \xi)]$

Corollary (Optimal Rate via Robust Online-to-Batch)

Suppose $\|\nabla f(\mathbf{w}, \xi)\| \leq G$ a.s. when $\xi \sim P$ is an inlier.

Then iterate average $\bar{\mathbf{w}}_T = \frac{1}{T} \sum_{t=1}^T \mathbf{w}_t$ of **OGD** + **Top-k Filter** achieves

$$\text{Risk}_P(\bar{\mathbf{w}}_T) - \min_{\mathbf{u} \in \mathcal{W}} \text{Risk}_P(\mathbf{u}) = O\left(DG\epsilon + DG\sqrt{\frac{\ln(1/\delta)}{T}}\right)$$

with P_ϵ -probability at least $1 - \delta$, for some k tuned for ϵ, δ, T .

Quantile Outliers

Which **extra assumptions** allow **sublinear** dependence on number of outliers k ?

- ▶ $\|g_t\| \leq L\|\mathbf{X}_t\|$ for i.i.d. \mathbf{X}_t (e.g. hinge loss, logistic loss)
- ▶ Inliers \mathcal{S}_p are rounds s.t. $\|\mathbf{X}_t\|$ less than p -quantile X_p

Quantile Outliers

Which **extra assumptions** allow **sublinear** dependence on number of outliers k ?

- ▶ $\|g_t\| \leq L\|\mathbf{X}_t\|$ for i.i.d. \mathbf{X}_t (e.g. hinge loss, logistic loss)
- ▶ Inliers \mathcal{S}_p are rounds s.t. $\|\mathbf{X}_t\|$ less than p -quantile X_p

Theorem (Sublinear Outlier Overhead)

Suppose ALG has regret bound $B_T(X)$, concave in T , if non-filtered \mathbf{X}_t have length at most X . Then **ALG** + **p -Quantile Filter** achieves

$$\mathbb{E} \left[\max_{\mathbf{u} \in \mathcal{W}} R_T(\mathbf{u}, \mathcal{S}_p) \right] \leq B_{pT}(X_p) + O \left(LDX_p \sqrt{p(1-p)T \ln T} + \ln(T)^2 \right).$$

p -Quantile Filter:

- ▶ Filter when $\|\mathbf{X}_t\| \geq$ lower-confidence bound on X_p

Summary

Robust regret: measure regret only on (unknown) inlier rounds

Price of Robustness = Overhead over usual regret rate:

- ▶ At most k adversarial outliers: $O(k)$
- ▶ p -Quantile outliers: $O(\sqrt{p(1-p)T \ln(T)} + \ln(T)^2)$

**PS. I am looking for a postdoc, starting anytime in 2022.
Please get in touch if you want to come to Amsterdam!**