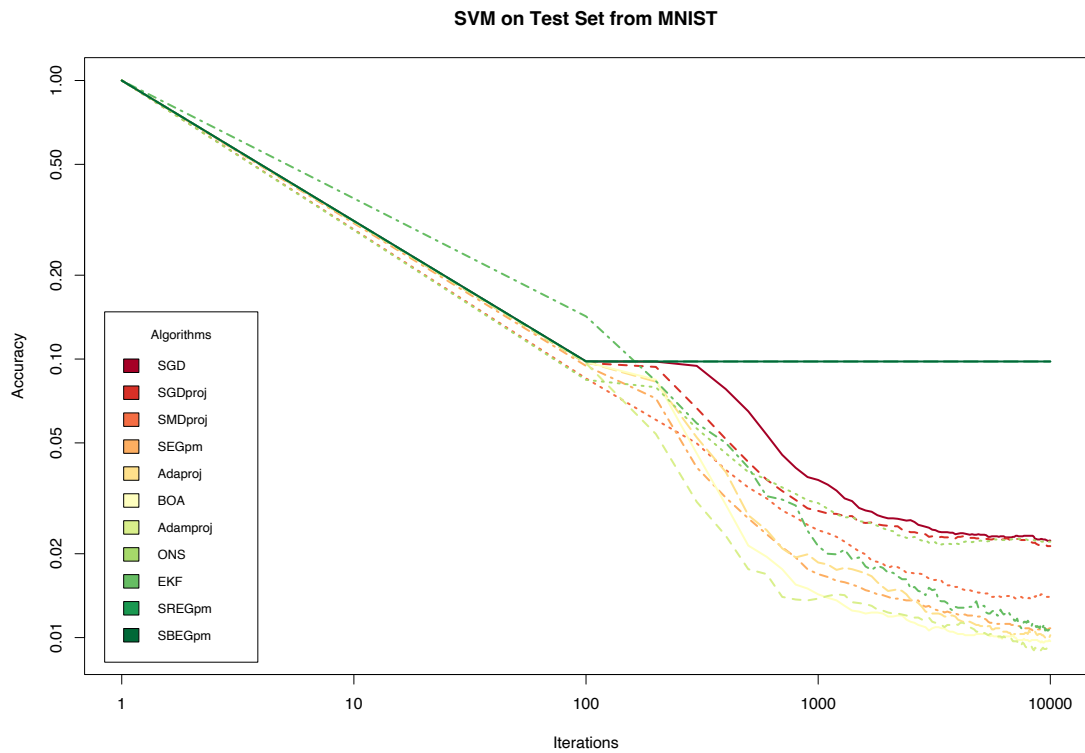

Online Convex Optimization

Olivier WINTENBERGER



Contents

I	Preliminaries	3
1	Basic concepts in Convex Optimization (CO)	5
1.1	Basic definition and setup	5
1.2	Gradient Descent algorithm (GD)	7
1.3	Applications	10
II	Online Convex Optimization	15
2	Online Gradient Descent for Online Convex Optimization (OCO)	17
2.1	The setting	17
2.2	Failure of Follow The Leader (FTL)	18
2.3	Online Gradient Descent (OGD)	19
2.4	Applications	21
3	Online Regularization	25
3.1	Online regularization	25
3.2	Online Mirror Descent	26
3.3	Specific OMD	29
III	Acceleration and exploration	45
4	Accelerated OCO algorithms	47
4.1	Momentum	47
4.2	Online Newton Step (ONS)	50
5	Exploration	59
5.1	Bandit Convex Optimization	59
5.2	Exp3 algorithm	61
5.3	Exp2 algorithm for OCO on $\mathcal{K} = B_1(z)$	64

These lectures notes are reproducing most of the 6 first Chapters of Hazan's book

Introduction to Online Convex Optimization

<https://sites.google.com/view/intro-oco/>

with small variations, especially in Section 5. Online Convex Optimization is the study of recursive algorithm and their theoretical guarantees called regret bounds. Due to the effectiveness of some algorithms of this vast class for training deep neural networks there is an excellent recent literature. Besides Hazan (2019), there is also the early Shalev-Shwartz et al. (2011) and the very recent Orabona (2019). Lattimore and Szepesvári (2020) is the reference for the treatment of bandit problems.

All the illustrations of these notes are maid on the MNIST handwritten digit dataset from

<http://yann.lecun.com/exdb/mnist/>

tuned into a classification problem recognizing the digit 0. The performances of linear SVM, trained on 60000 digits with different algorithms, are compared in terms of their accuracy on the test set of 10000 digits. The seed is fixed the same for all stochastic algorithms and all the experiments are ran on the R language from CRAN. It is available on

<http://wintenberger.fr/ens.html>

Part I
Preliminaries

Chapter 1

Basic concepts in Convex Optimization (CO)

In this chapter we fix some notation from the usual CO problem, cf *Convex Optimization*, Boyd et al. (2004), the more recent introductory notes Bubeck (2014) and *Remise à niveau. Calcul différentiel et optimisation*, the lecture notes of the course of Claire Boyer and Maxime Sangnier.

1.1 Basic definition and setup

We are interested in analyzing the performances of algorithms solving the CO problem. The key notion is convexity: Convexity of a set \mathcal{K}

$$\alpha x + (1 - \alpha)y \in \mathcal{K}, \quad x, y \in \mathcal{K}, \quad \alpha \in (0, 1),$$

and convexity of a function $f : \mathcal{K} \mapsto \mathbb{R}$

$$f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y).$$

In all the sequel \mathcal{K} will be assumed closed and bounded with diameter D

$$\|x - y\| \leq D, \quad x, y \in \mathcal{K}.$$

Here the norm $\|\cdot\| = \|\cdot\|_2$ is the Euclidean one over \mathbb{R}^d , $d \geq 1$, the other ℓ^p -norms being denoted $\|\cdot\|_p$. On a closed and bounded convex set a convex function admits a (non necessarily unique) minimum.

Definition 1 (CO problem). A CO problem (f, \mathcal{K}) is to approximate the minimum of f over \mathcal{K}

$$\min_{x \in \mathcal{K}} f(x),$$

or, alternatively, to approximate one of the minimizers

$$x^* \in \arg \min_{x \in \mathcal{K}} f(x) = \{x \in \mathcal{K}; f(x) = \min_{x \in \mathcal{K}} f(x)\}.$$

Another way of defining a CO problem is via its (sub-)gradient. A sub-gradient is a vector $\nabla f(x) \in \mathbb{R}^d$ satisfying the relation

$$f(y) \geq f(x) + \nabla f(x)^T (y - x), \quad x, y \in \text{Dom}(f), \quad (1.1)$$

where the domain $Dom(f) = \{x \in \mathbb{R}^d : f(x) < \infty\}$ of f is convex. For simplicity we assume that the sub-gradient is unique $\nabla f(x)$ at any point and $x \in \mathcal{K}$. We call $\nabla f(x)$ the gradient at $x \in \mathcal{K}$. If the minimizer is in the interior of \mathcal{K} then

$$x^* \in \arg \min_{x \in \mathcal{K}} f(x) \cap \overset{\circ}{\mathcal{K}} \iff \nabla f(x^*) = 0.$$

In generality they might be a problem on the boundary of \mathcal{K} .

Theorem 1 (simple Karush-Kuhn-Tucker (KKT)). *For every $y \in \mathcal{K}$ we have*

$$\nabla f(x^*)^T (y - x^*) \geq 0.$$

Proof. Assume that for some $y \in \mathcal{K}$ we have $\nabla f(x^*)^T (y - x^*) < 0$. Then consider $g(t) = f(x^* + t(y - x^*))$ so that $g'(0) = \nabla f(x^*)^T (y - x^*) < 0$. In particular for $t > 0$ sufficiently small we have $g(t) < g(0)$, thus $z = x^* + t(y - x^*) = ty + (1-t)x^* \in \mathcal{K}$ satisfies $f(z) < f(x^*)$ in contradiction with the definition of x^* . \square

It means that the gradient at the minimum points at the interior of the constrained set. We will use the projection $\Pi_{\mathcal{K}}(y) = \arg \min_{x \in \mathcal{K}} \|y - x\|$ the (convex) projection of y on \mathcal{K} . It is defined for the Euclidean norm and will be extended to other norms. As it is a CO problem, the projection is well defined but may be non explicit, see Grünewälder (2017)!

Exercise 1. *Show that $\Pi_{\mathcal{K}}(x)$ has an explicit solution $x/\|x\|$ if $\mathcal{K} = B_2(1) = B(1)$ the unitary Euclidian ball and $x \notin \mathcal{K}$.*

Theorem 2 (Pythagorean). *For any $z \in \mathcal{K}$ we have $\|y - z\|^2 \geq \|\Pi_{\mathcal{K}}(y) - z\|^2 + \|\Pi_{\mathcal{K}}(y) - y\|^2$.*

Proof. We have the CO problem $\Pi_{\mathcal{K}}(y) = \arg \min_{x \in \mathcal{K}} f(x)$ with $f(x) = \|x - y\|^2$ thus

$$\begin{aligned} \|y - z\|^2 - \|\Pi_{\mathcal{K}}(y) - z\|^2 &= \|y\|^2 - \|\Pi_{\mathcal{K}}(y)\|^2 + 2(\Pi_{\mathcal{K}}(y) - y)^T z \\ &= \|y\|^2 - \|\Pi_{\mathcal{K}}(y)\|^2 + \nabla f(\Pi_{\mathcal{K}}(y))^T z \\ &\geq \|y\|^2 - \|\Pi_{\mathcal{K}}(y)\|^2 + \nabla f(\Pi_{\mathcal{K}}(y))^T \Pi_{\mathcal{K}}(y) \\ &\geq \|y\|^2 - \|\Pi_{\mathcal{K}}(y)\|^2 + 2(\Pi_{\mathcal{K}}(y) - y)^T \Pi_{\mathcal{K}}(y) \\ &\geq \|y\|^2 + \|\Pi_{\mathcal{K}}(y)\|^2 - 2y^T \Pi_{\mathcal{K}}(y) \\ &\geq \|y - \Pi_{\mathcal{K}}(y)\|^2 \geq 0. \end{aligned}$$

We used the simple KKT theorem to derive the inequality. \square

Note that the above notions and theorem extends easily to any weighted quadratic forms

$$\|x\|_W^2 = x^T W x, \quad W \succ 0.$$

Here $W \succ 0$ means that W is a definite positive matrix of weights.

We assume the sub-gradients are bounded, there exists some $G > 0$ so that $\|\nabla f(x)\| \leq G$ for all $x \in \mathcal{K}$. Then f is Lipschitz (and continuous) $|f(y) - f(x)| \leq G\|y - x\|$, $x, y \in \mathcal{K}$.

We may assume more on the curvature of the function f , namely

- f is α -strongly convex if

$$f(y) \geq f(x) + \nabla f(x)^T (y - x) + \frac{\alpha}{2} \|y - x\|^2, \quad x, y \in Dom(f),$$

- f is β -smooth if

$$f(y) \leq f(x) + \nabla f(x)^T(y - x) + \frac{\beta}{2}\|y - x\|^2, \quad x, y \in \text{Dom}(f).$$

Exercise 2. Show that β -smoothness follows from ∇f is β -Lipschitz.

Remark 1. If the function f is twice differentiable we denote $\nabla^2 f(x)$ the Hessian $d \times d$ matrix at the point $x \in \mathcal{K}$ and

- f is α strongly convex iff $\nabla^2 f(x) \succeq \alpha I_d$ ($A \succeq 0$ meaning that A is a symmetric semi-definite positive matrix),
- f is β smooth iff $\nabla^2 f(x) \preceq \beta I_d$.

If f is α -strongly convex and β -smooth then f is γ -well-conditioned with $\gamma = \alpha/\beta \leq 1$.

A typical example is the quadratic loss $f(x) = \|x\|^2$ which is $\gamma = 1$ well-conditioned as $\alpha = \beta = 2$.

1.2 Gradient Descent algorithm (GD)

In view of (1.1), minimizing f from a given point $x \in \mathcal{K}$ is approximated by the CO problem on the *surrogate loss*, ie a simple approximation of the function $f(y) \approx f(x) + \nabla f(x)^T(y - x)$ in y . It is a linear function and one takes the step y from x in the opposite of the direction $x - \eta \nabla f(x)$ of the gradient so that $\nabla f(x)^T(y - x) < 0$. The role of η is to control the step-size, balancing between the gain in the surrogate CO problem (large η) and the quality of the approximation of the surrogate loss (small η).

Algorithm 1: Gradient Descent

Parameters: Epoch T , step-sizes (η_t) .

Initialization: Initial point $x_1 \in \mathcal{K}$.

For each iteration $t = 1, \dots, T$:

Iteration: Update

$$\begin{aligned} y_{t+1} &= x_t - \eta_t \nabla f(x_t), \\ x_{t+1} &= \Pi_{\mathcal{K}}(y_{t+1}). \end{aligned}$$

Return x_{T+1}

Let $h_t = f(x_t) - f(x^*)$, we have the rates

- γ -well-conditioned, $\eta_t = 1/\beta$, $h_T = O(e^{-\gamma T})$,
- β -smooth, $\eta_t = 1/\beta$, $h_T = O(\beta/T)$,
- α -strongly convex, $\eta_t = 1/(\alpha T)$, $h_T = O(1/(\alpha T))$,
- convex, $\eta_t = 1/\sqrt{T}$, $h_T = O(1/\sqrt{T})$.

The last two rates are optimal but the two first ones can be accelerated. Step sizes η_t depend on the curvature properties of the CO problem and the epoch T . The more the curvature of f is controlled the easier the CO problem (f, \mathcal{K}) .

Proof of the γ -well-conditioned unconstrained case. We show first the relation $h_t \leq h_{t-1}(1 - \alpha/\beta)$. By β -smoothness, we have

$$h_t - h_{t-1} = f(x_t) - f(x_{t-1}) \leq \nabla f(x_{t-1})^T(x_t - x_{t-1}) + \frac{\beta}{2}\|x_t - x_{t-1}\|^2.$$

We consider the upper-bound as a surrogate loss that we optimize in x_t . We obtain the desired gradient step $x_t = x_{t-1} - 1/\beta \nabla f(x_{t-1})$ and

$$\begin{aligned} h_t - h_{t-1} &\leq -\eta_t \|\nabla f(x_{t-1})\|^2 + \frac{\beta}{2} \eta_t^2 \|\nabla f(x_{t-1})\|^2 \\ &\leq -\frac{1}{2\beta} \|\nabla f(x_{t-1})\|^2. \end{aligned}$$

By strong convexity, we have

$$\begin{aligned} f(y) &\geq f(x) + \nabla f(x)^T(y - x) + \frac{\alpha}{2}\|x - y\|^2 \\ &\geq \min_{z \in \mathbb{R}^d} \{f(z) + \nabla f(x)^T(z - x) + \frac{\alpha}{2}\|x - z\|^2\} \\ &\geq f(x) - \frac{1}{2\alpha} \|\nabla f(x)\|^2 \end{aligned}$$

because $z^* = x - \frac{1}{\alpha} \nabla f(x)$. In particular taking $x = x_t$ and $y = x^*$ we get

$$\|\nabla f(x_t)\|^2 \geq 2\alpha(f(x_t) - f(x^*)) = 2\alpha h_t.$$

Combining both inequalities we obtain

$$\begin{aligned} h_t - h_{t-1} &\leq -\frac{1}{2\beta} \|\nabla f(x_{t-1})\|^2 \\ &\leq -\frac{\alpha}{\beta} h_{t-1}. \end{aligned}$$

Thus a recursive argument yields

$$h_T \leq h_{T-1}(1 - \alpha/\beta) \leq h_{T-1}e^{-\gamma} \leq \dots \leq h_1 e^{-\gamma(T-1)}$$

and the result follows. \square

Definition 2. *Regularizing the CO problem (f, \mathcal{K}) consists in adjoining a regularization function R strongly convex on \mathcal{K} and twice continuously differentiable so that $(f + R, \mathcal{K})$ becomes an easier CO problem.*

Consider the regularized problem $g(x) = f(x) + \alpha/2\|x - x_1\|^2$ when f is convex. Then g is α -strongly convex so that the CO problem gets easier and the error of the GD problem $h_T^g = g(x_T) - g(x_g^*)$ smaller. However the minimizer of the CO problem changes and we denote it x_g^* . Assuming that $x_1 \in \mathcal{K}$ we still have

$$\begin{aligned} f(x_T) - f(x^*) &= g(x_T) - g(x^*) + \alpha/2(\|x^* - x_1\|^2 - \|x_T - x_1\|^2) \\ &\leq g(x_T) - g(x_g^*) + \alpha/2D^2 && (g(x_g^*) \leq g(x^*)) \\ &\leq h_t^g + \alpha D^2/2. \end{aligned}$$

Exercise 3. *If f is convex and β -smooth show that $g = f + R$ with $R(x) = \alpha/2\|x - x_1\|^2$ is γ well conditioned. Deduce that $h_T = O(\beta \log T/T)$ choosing α carefully in the unconstrained CO problem.*

It is also possible to smooth the loss function thanks to randomization. Consider $\widehat{f}_\delta(x) = \mathbb{E}_{U \sim \mathcal{U}(B(1))}[f(x + \delta U)]$ where $\mathcal{U}(B(1))$ is the uniform distribution on the unit Euclidean ball. We have

Proposition 1. *The randomized version \widehat{f}_δ is dG/δ -smooth and a δG uniform approximation of f :*

$$|\widehat{f}_\delta(x) - f(x)| \leq \delta G, \quad x \in \mathcal{K}.$$

Proof. We use Stokes' theorem which is a multi-dimensional extension of the relation

$$\int_{-1}^1 f'(u) du = f(1) - f(-1) = (+1)f(+1) + (-1)f(-1).$$

Theorem 3 (Stokes' theorem). *For any continuously differentiable g we have*

$$\int_{B(1)} \nabla g(u) du = \int_{S(1)} g(v) v dv$$

where $B(1)$ and $S(1)$ are the unit Euclidean ball and sphere of \mathbb{R}^d , respectively.

Since $d|B(1)| = |S(1)|$ where $|\cdot|$ denotes the Lebesgue measure (in \mathbb{R}^d and \mathbb{R}^{d-1} respectively) we obtain

$$\int_{B(1)} \nabla f(x + \delta u) \frac{du}{|B(1)|} = \frac{d}{\delta} \int_{S(1)} f(x + \delta v) v \frac{dv}{|S(1)|}.$$

Then we have, interchanging derivative and expectation by domination and unicity of ∇f ,

$$\begin{aligned} \|\nabla \widehat{f}_\delta(x) - \nabla \widehat{f}_\delta(y)\| &= \|\mathbb{E}_{U \sim \mathcal{U}(B(1))}[\nabla f(x + \delta U)] - \mathbb{E}_{U \sim \mathcal{U}(B(1))}[\nabla f(y + \delta U)]\| \\ &= \frac{d}{\delta} \|\mathbb{E}_{V \sim \mathcal{U}(S(1))}[f(x + \delta V)V] - \mathbb{E}_{V \sim \mathcal{U}(S(1))}[f(y + \delta V)V]\| \\ &\leq \frac{d}{\delta} \mathbb{E}_{V \sim \mathcal{U}(S(1))}[\|(f(x + \delta V) - f(y + \delta V))V\|] \\ &\leq \frac{d}{\delta} \mathbb{E}_{V \sim \mathcal{U}(S(1))}[\|(f(x + \delta V) - f(y + \delta V))\| \|V\|] \\ &\leq \frac{d}{\delta} G \|x - y\| \mathbb{E}_{V \sim \mathcal{U}(S(1))}[\|V\|] \\ &\leq \frac{d}{\delta} G \|x - y\| \end{aligned}$$

and the β -smoothness follows. The approximation bound is easily computed using again Jensen's inequality and the Lipschitz property of f :

$$\begin{aligned} |\widehat{f}_\delta(x) - f(x)| &= |\mathbb{E}_{U \sim \mathcal{U}(B(1))}[f(x + \delta U) - f(x)]| \leq \mathbb{E}_{U \sim \mathcal{U}(B(1))}[|f(x + \delta U) - f(x)|] \\ &\leq G \mathbb{E}_{U \sim \mathcal{U}(B(1))}[\|\delta U\|] \leq \delta G. \end{aligned}$$

□

Consider the smoothed unconstrained problem $(\widehat{f}_\delta, \mathbb{R}^d)$. One deduces that

$$\begin{aligned} f(x_T) - f(x^*) &\leq \widehat{f}_\delta(x_T) - \widehat{f}_\delta(x^*) + 2\delta G \\ &\leq \widehat{f}_\delta(x_T) - \widehat{f}_\delta(x_{\widehat{f}_\delta}^*) + 2\delta G & (\widehat{f}_\delta(x_{\widehat{f}_\delta}^*) \leq \widehat{f}_\delta(x^*)) \\ &\leq h_t^{\widehat{f}_\delta} + 2\delta G. \end{aligned}$$

Exercise 4. *If f is α strongly convex show that \widehat{f}_δ is γ well conditioned. Deduce that $h_T = O(G^2 d \log T / (\alpha T))$ for δ well chosen in the unconstrained CO problem.*

Exercise 5. *Show that the rate of any CO problem is at most of the order $O(GD\sqrt{d \log T/T})$.*

1.3 Applications

1.3.1 Unconstrained CO problem

Consider the supervised classification problem of 2 classes $\{+1, 1\}$ and one observes labels b_i , $1 \leq i \leq n$, $b_i \in \{+1, 1\}$ together with explanatory variables $a_i \in \mathbb{R}^d$.

Examples:

- Natural Language Processing (NLP) for spam classification: a_i encodes the list of words in an email, d is the number of words in the language, $a_{i,j} = 1$ if the j word appears in the i th mail, $= 0$ else.
- MNIST: Handwritten digit database $n = 60000$ from a_i is a 28×28 grayscale image, $d = 784$ and one can consider two classes, 0 vs other digits ($b_i = 0$ if the digit is 0, else -1).

Definition 3. *Linear Support Vector Machine (SVM) are classifiers of the form $\text{sign}(x^T a_i)$ an hyperplane $x \in \mathbb{R}^d$.*

One wants to find the minimizer x of the accuracy

$$\mathbb{P}(\text{sign}(x^T a) \neq b) \approx x \mapsto \frac{1}{n'} \sum_{i=n+1}^{n+n'} \mathbb{1}_{\text{sign}(x^T a_i) \neq b_i}$$

over a test set $(a_i, b_i)_{n+1 \leq i \leq n+n'}$.

Because of the lack of convexity of the 0/1 loss, the hard-margin problem of optimizing

$$\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\text{sign}(x^T a_i) \neq b_i}$$

is non-polynomial. A common way of bypassing the issue is to relax the optimization problem to turn it to a CO problem.

Definition 4. *The hinge loss*

$$\ell_{a,b}(x) = \text{hinge}(bx^T a) = \max(0, 1 - bx^T a)$$

is a convex version of the 0 – 1 loss $\mathbb{1}_{\text{sign}(x^T a) \neq b} = \mathbb{1}_{bx^T a < 0}$.

Remark that the hinge loss is a convex function but not strongly convex with potentially multiple minimizers. We consider instead the regularized CO problem called the soft-margin problem

$$f(x) = \frac{1}{n} \sum_{i=1}^n \text{hinge}(bx^T a_i) + \frac{\lambda}{2} \|x\|^2.$$

It is a strongly convex CO.

1.3.2 ℓ^1 -ball constrained CO as dual of the LASSO

Consider $f(x) = \frac{1}{n} \sum_{i=1}^n \ell_i(x)$ to be minimized on the trained sample (ℓ_1, \dots, ℓ_n) assumed to be iid convex functions over \mathbb{R}^d . The aim is to minimize the unobserved risk $\mathbb{E}[\ell(x)]$. One can face generalization issues such as overfitting in high dimension.

Definition 5. *The information criteria (AIC, BIC) are penalized f of the form*

$$f(x) + \lambda \|x\|_0 = f(x) + \theta \sum_{i=1}^d \mathbb{1}_{x_i \neq 0}, \quad \theta > 0, x \in \mathbb{R}^d.$$

Due to the lack of convexity, it is a non-polynomial optimization problem. It is relaxed using the convex ℓ^1 -norm instead of $\|\cdot\|_0$.

Definition 6. *The LASSO problem is a penalized unconstrained CO problem of the form*

$$f(x) + \theta \|x\|_1, \quad \theta > 0, x \in \mathbb{R}^d.$$

LASSO is an unconstrained CO problem. Using a Lagrange dual argument one can turn it into a constrained CO problem.

Proposition 2. *If x^* is the minimizer of LASSO it is also the minimizer of*

$$\min_{\|x\|_1 \leq \tau} f(x), \quad (1.2)$$

$\tau = \|x^*\|_1$. Moreover if x^* is a minimizer of (1.2) then there exists θ^* such that it minimizes LASSO.

Proof. Assume x^* is the minimizer of the LASSO problem then for any $\|x\|_1 \leq \tau$ we have

$$\begin{aligned} f(x) &\geq f(x^*) + \theta(\|x^*\|_1 - \|x\|_1) \\ &\geq f(x^*) + \theta(\|x^*\|_1 - \tau) \\ &\geq f(x^*) \end{aligned}$$

when $\tau = \|x^*\|_1$. The second assertion requires the general KKT theorem

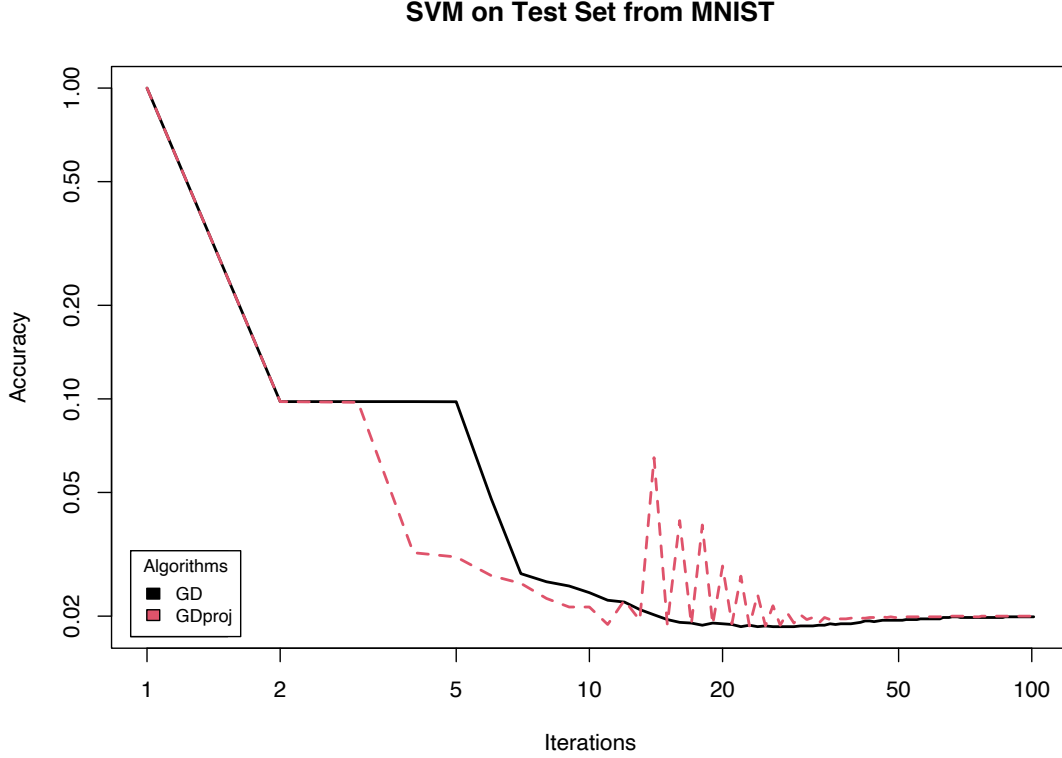
Theorem 4 (general KKT). *If (x^*, θ^*) is a saddle point of the Lagrangian $\mathcal{L}(x, \theta) = f(x) + \theta g(x)$ (minimum over $x \in \mathbb{R}^d$, maximum over $\theta \geq 0$) then x^* solves the CO problem $(f, \{x \in \mathbb{R}^d : g(x) \leq 0\})$. If g is convex, there exists x_0 such that $g(x_0) < 0$, then x^* solving the CO problem $(f, \{x \in \mathbb{R}^d : g(x) \leq 0\})$ is associated to θ^* such that (x^*, θ^*) is a saddle point of \mathcal{L} .*

Since $g(x) = \|x\|_1 - \tau$ is convex and $g(0) < 0$, a minimizer x^* of (1.2) is associated to θ^* such that (x^*, θ^*) is a saddle point of \mathcal{L} . In particular we have that x^* minimize $\mathcal{L}(\cdot, \theta^*)$, ie

$$f(x) + \theta(\|x\|_1 - \tau) \geq f(x^*) + \theta(\|x^*\|_1 - \tau), \quad x \in \mathbb{R}^d,$$

which is equivalent to x^* solving the LASSO problem. \square

Implementation of (projected) GD on MNIST with $\eta_t = 1/(\lambda t)$, regularization parameter $\lambda = 1/3$ and projection on $B_1(100) = \{x \in \mathbb{R}^d; \|x\|_1 \leq 100\}$, each iteration costs $O(nd + P)$ as it requires n gradients of dimension d and the projection on an ℓ^1 -ball of complexity P .



We notice that the accuracy are better for the projected version with a faster convergence rate but then their accuracies are deteriorating. It is due to overfitting and motivates early stopping methods.

1.3.3 Explicit projection on $B_1(z)$

We consider the CO problem $\Pi_{B_1(z)}(x) = \min_{y \in B_1(z)} \|y - x\|$. One cannot use GD since then a projection step is required... Luckily there exists an explicit solution. Consider the simpler projection on the simplex $\Pi_{\Lambda}(x)$ where $\Lambda = \{w \in \mathbb{R}_+^d; \sum_{i=1}^d w_i = 1\}$ for x such that $x_i \geq 0$ and $\|x\|_1 \geq 1$. We have the Lagrangian function

$$\mathcal{L}(w, \theta, \zeta) = \frac{1}{2} \|w - x\|^2 + \theta \left(\sum_{i=1}^d w_i - 1 \right) - \sum_{i=1}^d \zeta_i w_i$$

with parameters $w \in \mathbb{R}^d$, $\theta \in \mathbb{R}$ and $\zeta \in \mathbb{R}_+^d$. We compute its gradient

$$\nabla \mathcal{L}(w, \theta, \zeta) = \begin{pmatrix} w - x + \theta \mathbb{1} - \zeta \\ \sum_{i=1}^d w_i - 1 \\ -w \end{pmatrix}, \quad \mathbb{1} = (1, \dots, 1)^T.$$

Thus KKT provides

$$\begin{cases} w^* = x - \theta^* \mathbb{1} + \zeta^*, \\ \sum_{i=1}^d w_i^* = 1 \\ w_i^* = 0 \text{ or } w_i^* > 0 \text{ and } \zeta_i^* = 0. \end{cases}$$

To sum up we obtain the soft-thresholding $w_i^* = \text{SoftThreshold}(x_i, \theta^*) = \max(x_i - \theta^*, 0)$. Consider $g(\theta) = \sum_{i=1}^d (x_i - \theta)_+$. It is a non increasing piecewise linear function of θ with

range $[0, \|x\|_1]$ over the its domain $(0, \|x\|_1]$. The break points are $x_{(d)} \leq \dots \leq x_{(1)}$ the sorted coordinates with slopes $-d, \dots, -1$, respectively. Since $1 \in (0, \|x\|_1]$ there exists $d_0 \in \{1, \dots, d\}$ such that $g(x_{(j)}) < 1$ for all $j \leq d_0$ and $g(x_{(j)}) \geq 1$ for all $j > d_0$. Since $g(x_{(j)}) = \sum_{i=1}^{j-1} (x_{(i)} - x_{(j)})$ we obtain that

$$d_0 = \max \left\{ 1 \leq j \leq d : \sum_{i=1}^{j-1} (x_{(i)} - x_{(j)}) < 1 \right\}$$

is the number of non-thresholded coordinates in w^* . Moreover, since θ^* belongs to the interval $[x_{(d_0)}, x_{(d_0-1)})$ we have

$$g(\theta) = \sum_{i=1}^{d_0} (x_i - \theta)$$

and one deduces the expression of θ^* .

Algorithm 2: Projection On the Simplex Π_Λ

Input: $x \in \mathbb{R}^d$.

If $x \in \Lambda$

Then Return x .

Else

Sort $x_{(1)} \geq \dots \geq x_{(d)}$

Find $d_0 = \max \left\{ 1 \leq j \leq d : \sum_{i=1}^{j-1} (x_{(i)} - x_{(j)}) < 1 \right\}$

Define $\theta^* = \frac{1}{d_0} \left(\sum_{j=1}^{d_0} x_{(j)} - 1 \right)$

Return $w^* = \text{SoftThreshold}(x, \theta^*)$.

The projection over the ℓ^1 -ball follows easily from the one on the simplex by using symmetric arguments.

Algorithm 3: Projection On ℓ^1 -ball $\Pi_{B_1(z)}$

Input: $x \in \mathbb{R}^d$.

If $x \in B_1(z)$

Then Return x .

Else

$w^* = \Pi_\Lambda(|x|/z)$

Return $y = \text{sign}(x)w^*$.

Computational cost is $P = O(d \log(d))$ on average, as Quicksort.

Part II

Online Convex Optimization

Online Gradient Descent for Online Convex Optimization (OCO)

We extend the previous CO setting to the the OCO problem and analyses the online gradient descent.

2.1 The setting

We consider now a recursive setting. At each iteration t of the algorithm, the algorithm predicts x_t and then the loss function f_t is revealed, potentially varying through time. Then the algorithm incurs the loss $f_t(x_t)$ and its aim is to minimize its regret at any horizon T

$$\text{Regret}_T = \sum_{t=1}^T f_t(x_t) - \min_{x \in \mathcal{K}} \sum_{t=1}^T f_t(x),$$

its cumulative losses relative to the best strategy frozen through time.

Definition 7. *The full adversarial setting corresponds to f_t chosen by an adversary as the worst possible loss function given the past predictions x_t, x_{t-1}, \dots*

Example 1 (Rock, Paper, Scissor). *Consider the game with the following cost table where 0 denotes a draw, 1 denotes that the row player wins, and -1 denotes a column player victory:*

Algorithm \ Adversary	Rock	Paper	Scissor	= A
Rock	0	-1	1	
Paper	1	0	-1	
Scissor	-1	1	0	

where A is a 3×3 cost matrix. We consider then $\mathcal{K} = \{\text{Rock}, \text{Paper}, \text{Scissor}\}$, it is discrete (not convex). One randomizes the strategy by considering $x \in \Delta$, the simplex $\{x \in \mathbb{R}_+^3; x_1 + x_2 + x_3 = 1\}$ so that the strategy is $\mathbb{P}(\text{Rock}) = x_1 \dots$. Consider first the full adversarial setting; the algorithm have a randomized strategy x_t and plays Rock, Paper and Scissor $i_t = (0, 1, 2)$ according to the distribution x_t . Then the adversary chooses the worst move according to x_t (and not the sample move that she cannot predict). We obtain $f_t(x_t) =$

$\max_{y \in \Lambda} y^T A x_t$. It is a convex function of x_t since

$$\begin{aligned} \max_{y \in \Lambda} y^T A(\alpha x_1 + (1 - \alpha)x_2) &= \max_{y \in \Lambda} (\alpha y^T A x_1 + (1 - \alpha)y^T A x_2) \\ &\leq \alpha \max_{y \in \Lambda} y^T A x_1 + (1 - \alpha) \max_{y \in \Lambda} y^T A x_2. \end{aligned}$$

It is a full adversarial OCO called the zero-sum game.

Of course OCO also embeds much more gentle adversarial settings:

Proposition 3. *If $f_t = f$ is constant then, back to CO and the accuracy of the averaging $\bar{x}_T = \frac{1}{T} \sum_{t=1}^T x_t$ satisfies*

$$h_t^f = f(\bar{x}_T) - f(x^*) \leq \frac{1}{T} \sum_{t=1}^T f(x_t) - f(x^*) = \frac{\text{Regret}_T}{T}.$$

Definition 8. Stochastic OCO is if (f_t) is an independent random function sequence with constant mean called the risk $R = \mathbb{E}[f_1]$.

Proposition 4. *The accuracy of the averaging on the risk $R = \mathbb{E}[f_1]$ satisfies, on average,*

$$\mathbb{E}[h_T^R] \leq \mathbb{E}\left[\frac{\text{Regret}_T}{T}\right].$$

Proof. We have

$$\begin{aligned} \mathbb{E}[h_T^R] &= \mathbb{E}[R(\bar{x}_T) - R(x^*)] \leq \mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T R(x_t) - R(x^*)\right] \\ &\leq \mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T \mathbb{E}[f_t](x_t) - \mathbb{E}[f_t](x^*)\right] \\ &\leq \mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T \mathbb{E}[f_t(x_t) - f_t(x^*) \mid \mathcal{F}_{t-1}]\right] \\ &\leq \mathbb{E}\left[\frac{\text{Regret}_T}{T}\right], \end{aligned}$$

where one has to introduce the natural filtration $\mathcal{F}_t = \sigma(f_t, f_{t-1}, \dots, f_1)$ and uses the fact that x_t is \mathcal{F}_{t-1} -measurable. \square

The aim of OCO is to design algorithms with the best possible regret and at least sub-linear regrets

$$\text{Regret}_T = o(T).$$

2.2 Failure of Follow The Leader (FTL)

We call FTL the strategy from CO: at each t predicts

$$x_t = x_{t-1}^* \in \arg \min_{x \in \mathcal{K}} \sum_{k=1}^{t-1} f_k(x).$$

This strategy fails in the OCO setting. Consider $\mathcal{K} = [-1, 1]$, $f_1(x) = x/2$ and

$$f_k(x) = \begin{cases} -x & \text{if } k \text{ is even} \\ x & \text{else.} \end{cases}$$

Thus

$$\sum_{k=1}^{t-1} f_k(x) = \begin{cases} -x/2 & \text{if } t \text{ is odd} \\ x/2 & \text{else.} \end{cases}$$

so that FTL predicts $x_t = -1$ if t is odd and $= 1$ else and occurs $f_t(x_t) = 1$. Thus starting at $x_1 = 0$ the regret is

$$\text{Regret}_T = \sum_{t=1}^T f_t(x_t) - \min_{x \in \mathcal{K}} \sum_{t=1}^T f_t(x) = (T-1) + 1/2 = T - 1/2.$$

Note that the regret of the constant strategy $x_t = 0$ is $1/2$.

2.3 Online Gradient Descent (OGD)

This online version of GD has been introduced by Zinkevich (2003).

Algorithm 4: Online Gradient Descent

Parameters: Step-sizes (η_t) .

Initialization: Initial prediction $x_1 \in \mathcal{K}$.

For each recursion $t \geq 1$:

Predict: x_t

Incur: $f_t(x_t)$

Observe: $\nabla f_t(x_t)$

Recursion: Update

$$\begin{aligned} y_{t+1} &= x_t - \eta_t \nabla f_t(x_t), \\ x_{t+1} &= \Pi_{\mathcal{K}}(y_{t+1}). \end{aligned}$$

OGD succeeds where FTL fails:

Theorem 5. *OGD with $\eta_t = \frac{D}{G\sqrt{t}}$ satisfies*

$$\text{Regret}_T \leq \frac{3}{2}GD\sqrt{T}$$

Proof. We start with the gradient trick $f_t(x_t) - f_t(x^*) \leq \nabla f_t(x_t)(x_t - x^*)$, $t \geq 1$. Thus we will estimate the linear regret

$$\sum_{t=1}^T \nabla f_t(x_t)(x_t - x^*)$$

Using the updates, we have

$$\begin{aligned} \|x_{t+1} - x^*\|^2 &\leq \|\Pi_{\mathcal{K}}(x_t - \eta_t \nabla f_t(x_t)) - x^*\|^2 \\ &\leq \|x_t - \eta_t \nabla f_t(x_t) - x^*\|^2 \\ &\leq \|x_t - x^*\|^2 + \eta_t^2 \|\nabla f_t(x_t)\|^2 - 2\eta_t \nabla f_t(x_t)^T (x_t - x^*). \end{aligned}$$

We get, whatever is x^* ,

$$2\eta_t \nabla f_t(x_t)^T(x_t - x^*) \leq \|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2 + \eta_t^2 G^2. \quad (2.1)$$

One deduces that ($1/\eta_0 = 0$ by convention and $\|x_{t+1} - x^*\|^2 \geq 0$)

$$\begin{aligned} 2 \sum_{t=1}^T \nabla f_t(x_t)^T(x_t - x^*) &\leq \sum_{t=1}^T \frac{\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2}{\eta_t} + G^2 \sum_{t=1}^T \eta_t \\ &\leq \sum_{t=1}^T \|x_t - x^*\|^2 \left(\frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) + G^2 \sum_{t=1}^T \eta_t \\ &\leq D^2 \sum_{t=1}^T \left(\frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) + G^2 \sum_{t=1}^T \eta_t \\ &\leq \frac{D^2}{\eta_T} + G^2 \sum_{t=1}^T \eta_t \\ &\leq 3DG\sqrt{T}. \end{aligned}$$

We use that $\sum_{t=1}^T \eta_t \leq 2\sqrt{T}$. □

Remark 2. Note that if η is constant then a similar argument yields the upper bound

$$\frac{1}{2} \left(\frac{\|x_1 - x^*\|^2}{\eta} + \eta \sum_{t=1}^T \|\nabla f_t(x_t)\|^2 \right)$$

which is minimized for

$$\eta = \frac{\|x_1 - x^*\|}{\sqrt{\sum_{t=1}^T \|\nabla f_t(x_t)\|^2}}$$

and get the optimal regret bound

$$D \sqrt{\sum_{t=1}^T \|\nabla f_t(x_t)\|^2}.$$

However this bound is not achievable since the learning rate η is tuned knowing the gradients $\nabla f_t(x_t)$, $1 \leq t \leq T$, which are not observed.

Exercise 6. Compute the regret in the OCO where FTL fails. Interpret.

Moreover in favorable cases one can accelerate OGD:

Definition 9 (Strongly convex OCO). We consider the OCO problem over \mathcal{K} and we assume the existence of $\alpha > 0$ so that the f_t are α -strongly convex.

OGD satisfies an optimal regret bound $O(\log T)$ by modifying the step-sizes (learning rates) accordingly:

Theorem 6. Assume the strongly convex OCO problem, then OGD with step sizes $\eta_t = 1/(\alpha t)$ satisfies

$$\text{Regret}_t \leq \frac{G^2}{2\alpha} (1 + \log T).$$

Proof. We write

$$\text{Regret}_T = \sum_{t=1}^T f_t(x_t) - f_t(x^*)$$

so that by α -strong convexity we get the improved gradient trick

$$2(f_t(x_t) - f_t(x^*)) \leq 2\nabla f_t(x_t)^T(x_t - x^*) - \alpha\|x_t - x^*\|^2.$$

Using (2.1), namely

$$2\eta_t \nabla f_t(x_t)^T(x_t - x^*) \leq \|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2 + \eta_t^2 G^2$$

we get

$$\begin{aligned} 2(f_t(x_t) - f_t(x_T^*)) &\leq \frac{\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2}{\eta_t} - \alpha\|x_t - x^*\|^2 + \eta_t G^2 \\ &\leq \alpha(t-1)\|x_t - x^*\|^2 - \alpha t\|x_{t+1} - x^*\|^2 + \frac{G^2}{\alpha t}. \end{aligned}$$

Thus one gets a telescoping argument when bounding the regret

$$\begin{aligned} 2\text{Regret}_T &= \sum_{t=1}^T 2(f_t(x_t) - f_t(x^*)) \\ &\leq \sum_{t=1}^T \alpha(t-1)\|x_t - x^*\|^2 - \alpha t\|x_{t+1} - x^*\|^2 + \frac{G^2}{\alpha t} \\ &\leq -\alpha T\|x_{T+1} - x^*\|^2 + \frac{G^2}{\alpha}(1 + \log T) \end{aligned}$$

since $\sum_{t=1}^T t^{-1} \leq 1 + \log T$. The desired result follows. \square

Note that the regret bounds for OGD in the convex and strongly convex OCO problem are optimal.

2.4 Applications

2.4.1 Stochastic Gradient Descent (SGD)

Consider the CO problem (f, \mathcal{K}) . Instead of using ∇f we use a noisy version of the gradient $\widehat{\nabla} f$ so that $\mathbb{E}[\widehat{\nabla} f(x)] = \nabla f(x)$ and $\mathbb{E}[\|\widehat{\nabla} f(x)\|^2] \leq G^2$, independent of anything else. The approximation $\widehat{\nabla} f$ is unbiased with bounded variance and the setting is called Stochastic Optimisation (SO).

Example 2. Consider the unconstrained CO problem (f, \mathbb{R}^d) and the smoothed version $\hat{f}_\delta(x) = \mathbb{E}[f(x + \delta U)]$. Then $(d/\delta)f(x + \delta V)V$ is an unbiased estimator of $\nabla \hat{f}_\delta(x)$ due to Stokes' theorem. Remark we have $\mathbb{E}[\|\nabla \hat{f}_\delta\|^2] \leq (d/\delta)^2 \mathbb{E}[\|f(x + \delta V)V\|^2] = O((dG)^2)$.

Example 3. Consider the CO problem (f, \mathcal{K}) where $f(x) = \frac{1}{n} \sum_{i=1}^n \ell_i(x)$ as in the SVM classification. Then each step of a GD costs $O(nd)$ since it requires the query of n gradients $\nabla \ell_i$. Instead sample randomly uniformly $I \in \{1, \dots, n\}$ and use $\widehat{\nabla} f = \nabla \ell_I$. We have

$$\mathbb{E}[\widehat{\nabla} f(x)] = \sum_{i=1}^n \nabla \ell_i(x) \mathbb{P}(i = n) = \frac{1}{n} \sum_{i=1}^n \nabla \ell_i(x) = \nabla f(x)$$

and

$$\mathbb{E}[\|\widehat{\nabla}f(x)\|^2] = \sum_{i=1}^n \|\nabla\ell_i(x)\|^2 \mathbb{P}(i = n) = \frac{1}{n} \sum_{i=1}^n \|\nabla\ell_i(x)\|^2 \leq G^2$$

as soon as $\|\nabla\ell_i(x)\| \leq G$. Each step of a Stochastic GD (SGD) on $\nabla\ell_I$ costs $O(d)$.

Remark that by Jensen's inequality we also have $\|\nabla f\|^2 \leq \mathbb{E}[\|\widehat{\nabla}f(x)\|^2]$ in the example above.

Proposition 5. Any SO problem using iid unbiased approximations $\widehat{\nabla}f_t$ at each round t reduces to a stochastic OCO problem by considering $\nabla f_t(x_t) = \widehat{\nabla}f_t(x_t)$.

Proof. A SO problem requires that the approximations $\widehat{\nabla}f_t$ are all unbiased and independent and the optimizer introduces the randomness and chooses the distribution of $\widehat{\nabla}f_t$. In the stochastic OCO problem it is the nature which is random and chooses the distribution of f_t with mean $R = \mathbb{E}[f_1]$ and independent. Forgetting that the distribution is chosen by the optimizer, an algorithm robust to any choice of distribution will have good accuracy on the risk R in the stochastic OCO problem. It implies the same accuracy bound on the deterministic function $f = R$ whatever the optimizer chooses as $\widehat{\nabla}f_t$. \square

Based on this equivalence and on Proposition 4, SGD is a stochastic gradient descent together with an averaging step.

Algorithm 5: Stochastic Gradient Descent, Robbins and Monro (1951)

Parameters: Epoch T , step-sizes (η_t) .

Initialization: Initial point $x_1 \in \mathcal{K}$.

For each iteration $t = 1, \dots, T$:

Iteration: Sample $\widehat{\nabla}f_t$ independently of the rest.

Update

$$y_{t+1} = x_t - \eta_t \widehat{\nabla}f_t(x_t),$$

$$x_{t+1} = \Pi_{\mathcal{K}}(y_{t+1}).$$

Return: $\bar{x}_{T+1} = \frac{1}{T+1} \sum_{t=1}^{T+1} x_t$

The SGD is an iterative algorithm that can be studied via the stochastic OCO setting.

Theorem 7. SGD algorithm applied to the CO problem (f, \mathcal{K}) with $\eta_t = D/(G\sqrt{t})$ have an accuracy satisfying, on average,

$$\mathbb{E}[h_T^R] \leq \frac{3GD}{2\sqrt{T}}.$$

SGD algorithm applied to α -strongly CO problem (f, \mathcal{K}) with $\eta_t = 1/(\alpha t)$ have an accuracy satisfying, on average,

$$\mathbb{E}[h_T^R] \leq \frac{G^2}{2\alpha} \frac{1 + \log T}{T}.$$

The original CO setting is deterministic, the randomness comes from the sample of the approximations $\widehat{\nabla}f$ at each iteration. The expectation holds on this randomness.

The results on SGD follows from the regret bounds for OCO together with the online to batch conversion provided in Proposition 4.

The bounds are optimal, up to a log term. We gain in term of complexity; assume we are interested in an average accuracy of order $\epsilon > 0$ in the strongly convex CO problem (f, \mathcal{K}) . GD and SGD would require both $T = O(\epsilon^{-1})$ iterations. However the cost of each iteration is $O(nd)$ and $O(d)$ for GD and SGD, respectively, ending to a total cost of $O(nd\epsilon^{-1})$ and $O(d\epsilon^{-1})$, respectively. When n is large SGD is much more efficient on average!

2.4.2 Soft margin problem for linear SVM

Recall the soft margin problem which is a strongly convex CO problem with

$$f = \frac{1}{n} \sum_{i=1}^n \ell_{a_i, b_i}(x) + \frac{\lambda}{2} \|x\|^2$$

where $\ell_{a,b}(x) = \max(0, 1 - bx^T a)$.

Sampling I uniformly over $\{1, \dots, n\}$ one gets the approximation

$$\widehat{\nabla f_I}(x) = \nabla \ell_{a_I, b_I}(x) + \lambda x$$

which is unbiased. Since the learning rate is tuned as $1/(\lambda t)$ we get the SGD for solving the soft margin problem

Algorithm 6: SGD for linear SVM.

Parameters: Epoch T , radius $z > 0$, regularization parameter $\lambda > 0$

Initialization: Initial point $x_1 = 0$.

Sample uniformly iid: $(I_t)_{1 \leq t \leq T}$ from $\{1 \leq i \leq n\}$

For each iteration $t = 1, \dots, T$:

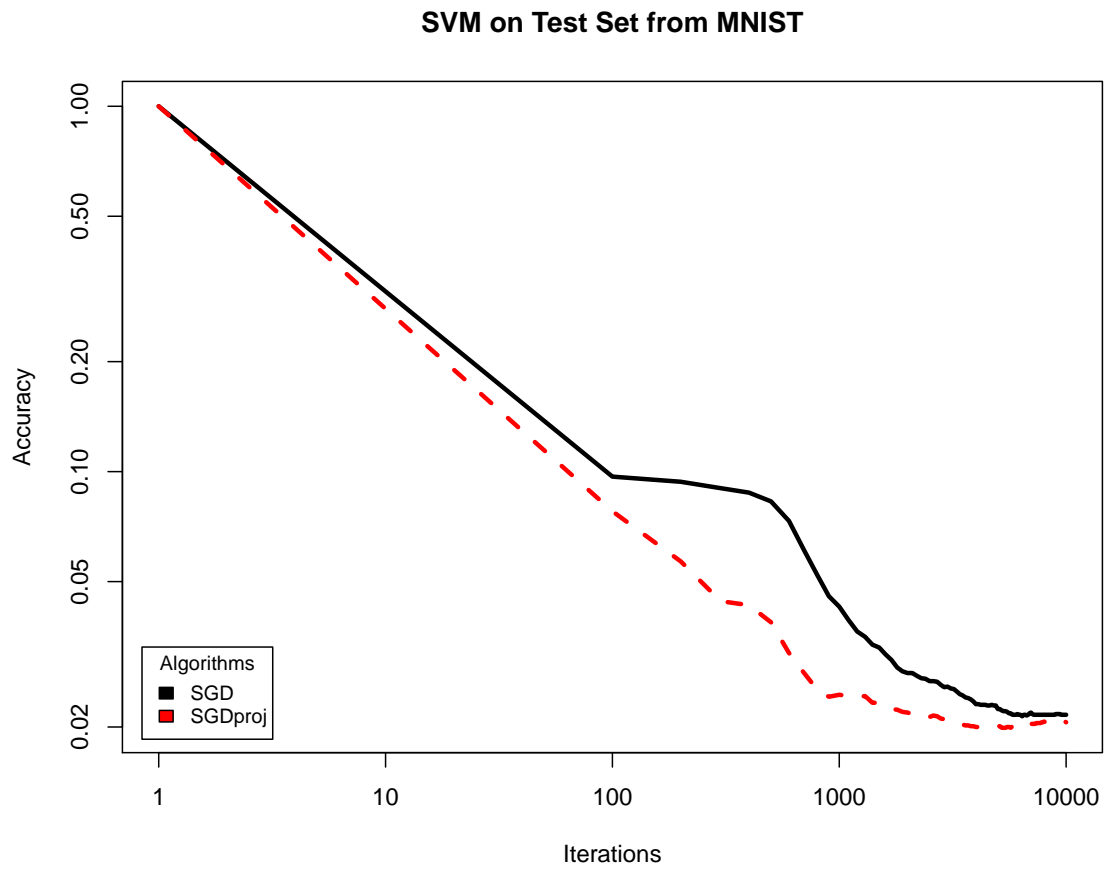
Iteration: Update

$$y_{t+1} = (1 - 1/t)x_t - \frac{\nabla \ell_{a_{I_t}, b_{I_t}}(x_t)}{\lambda t},$$

$$x_{t+1} = \Pi_{B_1(z)}(y_{t+1}).$$

Return: $\bar{x}_{T+1} = \frac{1}{T+1} \sum_{t=1}^{T+1} x_t$

The accuracy of the algorithm is on average $O(1/T)$ neglecting log terms. Implementation of (projected) SGD on MNIST with regularization parameter $\lambda = 1/3$ and projection on $B_1(100) = \{x \in \mathbb{R}^d; \|x\|_1 \leq 100\}$, each iteration costs $O(d)$ and the relative speed 1/1000 compared to GD. The



Online Regularization

3.1 Online regularization

We develop a general strategy for designing efficient OCO algorithms. The basic idea is to regularize FTL online so that it does not change too abruptly. Variants of OGD that fall in this class are Regularized FTL OCO algorithms.

Let R be a strongly convex regularization function twice continuously differentiable on \mathcal{K} . We replace the instance of FTL

$$x_t^* = \arg \min_{x \in \mathcal{K}} \sum_{k=1}^t f_k(x)$$

with x_{t+1} such as

$$x_{t+1} = \arg \min_{x \in \mathcal{K}} \left\{ \sum_{k=1}^t \nabla f_k(x_k)^T x + \frac{1}{\eta} R(x) \right\}.$$

The explanation consists in two steps; the first one is to change the cumulative loss up to t with a surrogate loss thanks to the gradient trick:

$$\sum_{k=1}^t f_k(x) - \sum_{k=1}^t f_k(x^*) \leq \sum_{k=1}^t \nabla f_k(x)^T (x - x^*).$$

and replacing the unobserved $\sum_{k=1}^t \nabla f_k(x)$ with approximations $\sum_{k=1}^t \nabla f_k(x_k)$. The obtained surrogate loss is linear

$$\sum_{k=1}^t \nabla f_k(x_k)^T (x - x^*).$$

The second step consists in regularizing this simple linear (convex) loss

$$\sum_{k=1}^t \nabla f_k(x_k)^T (x - x^*) + \frac{1}{\eta} R(x).$$

Doing so, we aim at obtaining an explicit formula for RFTL x_{t+1} (use of a simple surrogate

loss) and a more stable (regularized) version of FTL. We obtain

Algorithm 7: Regularized Follow The Leader, Shalev-Shwartz and Singer (2007)

Parameters: Regularization function R , step-size $\eta > 0$.

Initialization: Initial prediction $x_1 \in \mathcal{K}$.

For each recursion $t \geq 1$:

Predict: x_t

Incur: $f_t(x_t)$

Observe: $\nabla f_t(x_t)$

Recursion: Update

$$x_{t+1} = \arg \min_{x \in \mathcal{K}} \left\{ \sum_{k=1}^t \nabla f_k(x_k)^T x + \frac{1}{\eta} R(x) \right\}.$$

RFTL is a class of OCO algorithms. One has to specify R to specify the properties of RFTL.

3.2 Online Mirror Descent

OMD is an alternative way of defining RFTL in a more explicit way. For that we use the convex duality defined as

Definition 10. Let R be a regularization function defined on the convex set \mathcal{K} then its convex conjugate R^* is defined on the dual space $\mathcal{K}^* = \{\nabla R(x), x \in \mathcal{K}\}$ as

$$R^*(x^*) = \max_{x \in \mathcal{K}} \{x^T x^* - R(x)\}, \quad x^* \in \mathcal{K}^*.$$

Exercise 7. Prove that R^* is convex, $R^*(x^*) + R(y) \geq y^T x^*$ (the Fenchel-Young inequality) and

$$\nabla R^*(x^*) = \arg \max_{x \in \mathcal{K}} \{x^T x^* - R(x)\}.$$

Compute the conjugate of $R(x) = |x|^p/p$ for any $1 < p < \infty$ in the unconstrained case.

It is very natural to introduce the duality since the control of the scalar product provided in the Fenchel-Young inequality yields a regret bound from the gradient trick

$$\sum_{t=1}^T \nabla f_t(x_t)^T x \leq R^* \left(\sum_{t=1}^T \nabla f_t(x_t) \right) + R(x).$$

The OMD algorithm is designed for obtaining a good bound over $R^* \left(\sum_{t=1}^T \nabla f_t(x_t) \right)$.

OMD is an Online Gradient Descent in the convex "dual" space through the regularization function R defined as $\{\nabla R(x), x \in \text{Dom}(\nabla R)\}$, the space of the gradients of R (with restrictions on the domain of definition of the gradients so that they exists). The projection back to the primal space \mathcal{K} is driven by the Bregman divergence of R rather than by the usual Euclidian norm.

Definition 11. The Bregman divergence associated to the regularization function R is defined as

$$B_R(y||x) = R(y) - R(x) - \nabla R(x)^T (y - x).$$

The Bregman divergence shares some similarities with weighted Euclidian norms

$$\|x\|_W^2 = x^T W x, \quad W \succ 0.$$

Exercise 8. Show that $B_R(x||y) \geq 0$ and $B_R(x||y) = 0$ iff $x = y$.

Show that if R is twice continuously differentiable then

$$B_R(x||y) = \frac{1}{2} \|x - y\|_z^2$$

where $\|\cdot\|_z$ is some local norm

$$\|x - y\|_z^2 = (x - y)^T \nabla^2 R(z) (x - y)$$

for $z \in \mathcal{K}$ on the segment $[x, y]$.

Thus OMD will be specified through the properties of the Bregman divergence of R

Algorithm 8: Online Mirror Descent (lazy version), Hazan and Kale (2010)

Parameters: Regularization function R , step-size $\eta > 0$.

Initialization: Initial prediction $x_1 = \arg \min_{x \in \mathcal{K}} B_R(x||y_1)$ with $y_1 \in \mathbb{R}^d$ such that $\nabla R(y_1) = 0$.

For each recursion $t \geq 1$:

Predict: x_t

Incur: $f_t(x_t)$

Observe: $\nabla f_t(x_t)$

Recursion: Update

$$\begin{aligned} \nabla R(y_{t+1}) &= \nabla R(y_t) - \eta \nabla f_t(x_t), \\ x_{t+1} &= \arg \min_{x \in \mathcal{K}} B_R(x||y_{t+1}). \end{aligned}$$

Note that the gradient step makes a move over the dual space of the gradients $\{\nabla R(x), x \in \text{Dom}(\nabla R)\}$ which should be equal to \mathbb{R}^d so that y_{t+1} is defined for any potential value of $\eta \nabla f_t(x_t)$. One way to ensure that is to consider R so that $\lim_{x \rightarrow \partial \text{Dom}(\nabla R)} \|\nabla R(x)\|^2 = \infty$ for any point in the frontier of the domain $\partial \text{Dom}(\nabla R)$. Such regularization functions are called *Legendre* regularization function and are implicitly considered in the following.

Theorem 8. The OMD (lazy version) is equivalent to RFTL.

Proof. We prove that

$$\arg \min_{x \in \mathcal{K}} B_R(x||y_t) = \arg \min_{x \in \mathcal{K}} \left\{ \sum_{k=1}^t \nabla f_k(x_k)^T x + \frac{1}{\eta} R(x) \right\}.$$

Observe first that by recursion

$$\nabla R(y_t) = \nabla R(y_{t-1}) - \eta \nabla f_{t-1}(x_{t-1}) = -\eta \sum_{k=1}^{t-1} \nabla f_k(x_k).$$

Hence

$$\begin{aligned}
B_R(x||y_t) &= R(x) - R(y_t) - \nabla R(y_t)^T(x - y_t) \\
&= R(x) - R(y_t) + \eta \sum_{k=1}^{t-1} \nabla f_k(x_k)^T(x - y_t) \\
&= \eta \sum_{k=1}^{t-1} \nabla f_k(x_k)^T x + R(x) - R(y_t) - \underbrace{\eta \sum_{k=1}^{t-1} \nabla f_k(x_k)^T y_t}_{\text{independent of } x}
\end{aligned}$$

The desired results follows. \square

Denoting $\theta_t = \nabla R(y_t)$ and using the relation

$$\nabla R^*(x) = \arg \max_{y \in \mathcal{K}} \{y^T x - R(y)\}.$$

we have the equivalent formulation of OMD: We update

$$\theta_{t+1} = \theta_t - \eta \nabla f_t(x_t), \quad x_{t+1} = \nabla R^*(\theta_{t+1}),$$

from the initialization $\theta_1 = 0$ and $x_1 = \arg \max_{y \in \mathcal{K}} \{y^T \theta_{t+1} - R(y)\} = \nabla R^*(\theta_0)$. From that simple formulation and using similar argument than for proving the OGD regret bound, we get

Theorem 9. *OMD (lazy version) and thus RFTL satisfy the regret bound, of any $u \in \mathcal{K}$,*

$$\sum_{t=1}^T f_t(x_t) - \sum_{t=1}^T f_t(u) \leq \frac{\eta}{2} \sum_{t=1}^T \|\nabla f_t(x_t)\|_t^{*2} + \frac{R(u) - R(x_1)}{\eta},$$

where $\|\cdot\|_t^{*2} = \|\cdot\|_{\nabla^2 R^*(z_t^*)}^2$ for R^* the convex conjugate of R and z_t^* some point in \mathcal{K}^* .

Proof. We use the gradient trick

$$\sum_{t=1}^T f_t(x_t) - \sum_{t=1}^T f_t(u) \leq \sum_{t=1}^T \nabla f_t(x_t)^T(x_t - u).$$

We introduce the mirror analysis by introducing the dual θ_t :

$$\begin{aligned}
\sum_{t=1}^T \nabla f_t(x_t)^T(x_t - u) &= -\frac{1}{\eta} \sum_{t=1}^T (\theta_{t+1} - \theta_t)^T x_t + \frac{1}{\eta} \theta_{T+1}^T u \\
&= -\frac{1}{\eta} \sum_{t=1}^T (\theta_{t+1} - \theta_t)^T \nabla R^*(\theta_t) + \frac{R(u) + R^*(\theta_{T+1})}{\eta}
\end{aligned}$$

as $x_t = \nabla R^*(\theta_t)$ and applying Young's inequality. By definition of the Bregman divergence

$$\begin{aligned}
\sum_{t=1}^T \nabla f_t(x_t)^T(x_t - u) &= \frac{1}{\eta} \sum_{t=1}^T (R^*(\theta_t) - R^*(\theta_{t+1}) + B_{R^*}(\theta_{t+1}||\theta_t)) + \frac{R(u) + R^*(\theta_{T+1})}{\eta} \\
&= \frac{1}{\eta} \sum_{t=1}^T B_{R^*}(\theta_{t+1}||\theta_t) + \frac{R(u) + R^*(\theta_1)}{\eta}.
\end{aligned}$$

One recognizes $B_{R^*}(\theta_{t+1}||\theta_t) = \|\theta_{t+1} - \theta_t\|_t^{*2}/2 = \eta^2 \|\nabla f_t(x_t)\|_t^{*2}/2$ for z_t^* on the segment $[\theta_t, \theta_{t+1}]$ and $R^*(\theta_1) = \max_{y \in \mathcal{K}} y^T \theta_1 - R(y) = \max_{y \in \mathcal{K}} -R(y) = -R(x_1)$. The desired result follows. \square

3.3 Specific OMD

3.3.1 Quadratic Regularization

Online Mirror Descent is usually thought as RFLT associated to quadratic R . Consider $R(x) = \frac{1}{2}\|x - x_1\|^2$ for an arbitrary $x_1 \in \mathcal{K}$ and $\eta > 0$. Then $\nabla R(y_1) = (y_1 - x_1) = 0$ iff $y_1 = x_1$. Moreover

$$\begin{aligned} B_R(x|y) &= \frac{1}{2}\|x - x_1\|^2 - \frac{1}{2}\|y - x_1\|^2 - (y - x_1)^T(x - y) \\ &= \frac{1}{2}\|x - y + y - x_1\|^2 - \frac{1}{2}\|y - x_1\|^2 - (y - x_1)^T(x - y) \\ &= \frac{1}{2}\|x - y\|^2 \end{aligned}$$

so that

$$x_{t+1} = \arg \min_{x \in \mathcal{K}} B_R(x|y_{t+1}) = \Pi_{\mathcal{K}}(y_{t+1}).$$

We have $\nabla R(y_{t+1}) = (y_{t+1} - x_1) = \theta_{t+1} = -\eta \sum_{k=1}^t \nabla f_k(x_k)$ so that

$$y_{t+1} = y_t - \eta \nabla f_t(x_t), \quad y_1 = x_1.$$

Thus OMD for quadratic regularization function is an unconstrained OGD then projected on \mathcal{K} at each iteration.

Algorithm 9: Online Mirror Descent (for quadratic R)

Parameters: step-size $\eta > 0$.

Initialization: Initial prediction $x_1 = y_1 \in \mathcal{K}$.

For each recursion $t \geq 1$:

Predict: x_t

Incur: $f_t(x_t)$

Observe: $\nabla f_t(x_t)$

Recursion: Update

$$\begin{aligned} y_{t+1} &= y_t - \eta \nabla f_t(x_t), \\ x_{t+1} &= \Pi_{\mathcal{K}}(y_{t+1}). \end{aligned}$$

Remark 3. An agile version of the general OMD consists in replacing the recursion step $\nabla R(y_{t+1}) = \nabla R(y_t) - \eta \nabla f_t(x_t)$ with $\nabla R(y_{t+1}) = \nabla R(x_t) - \eta \nabla f_t(x_t)$. It moves faster than the lazy version of the OMD for the same learning rates. Note that the agile version of OMD with quadratic regularizer R is equivalent to OGD for the same learning rate.

Exercise 9. Show that in the unconstrained CO problem the lazy and agile versions of the OMD coincide. On the contrary, imagine a CO problem such that both OMD are projected at each step over the unit Euclidian ball $\mathcal{K} = B_1$. Then show that the lazy version with fixed $\eta = 1$ coincides with the agile version only if its learning rate at time t is $\|\sum_{s=1}^{t-1} \nabla f_s(x_s)\|^{-1} < 1$.

Since $R^*(x^*) = \frac{1}{2}(\|x^* + x_1\|^2 - \|x_1\|^2) \leq \frac{1}{2}\|x^*\|^2$ for any $x^* \in \mathcal{K}^*$ so that $x^* = \nabla R(x) =$

$x - x_1$ for some $x \in \mathcal{K}$, we get the regret bound

$$\begin{aligned} \sum_{t=1}^T f_t(x_t) - \sum_{t=1}^T f_t(u) &\leq \frac{\eta}{2} \sum_{t=1}^T \|\nabla f_t(x_t)\|^2 + \frac{\|u - x_1\|^2}{2\eta} \\ &\leq \frac{1}{2} \left(\eta T G^2 + \frac{D^2}{\eta} \right) \\ &\leq GD\sqrt{T} \end{aligned}$$

choosing $\eta = D/(G\sqrt{T})$. The lazy and agile (OGD) versions of the OMD got similar regret bound despite different learning rates.

Algorithm 10: SMD for linear SVM

Parameters: Epoch T , radius $z > 0$

Initialization: Initial point $x_1 = y_1 = 0$.

Sample uniformly iid: $(I_t)_{1 \leq t \leq T}$ from $\{1 \leq i \leq n\}$

For each iteration $t = 1, \dots, T$:

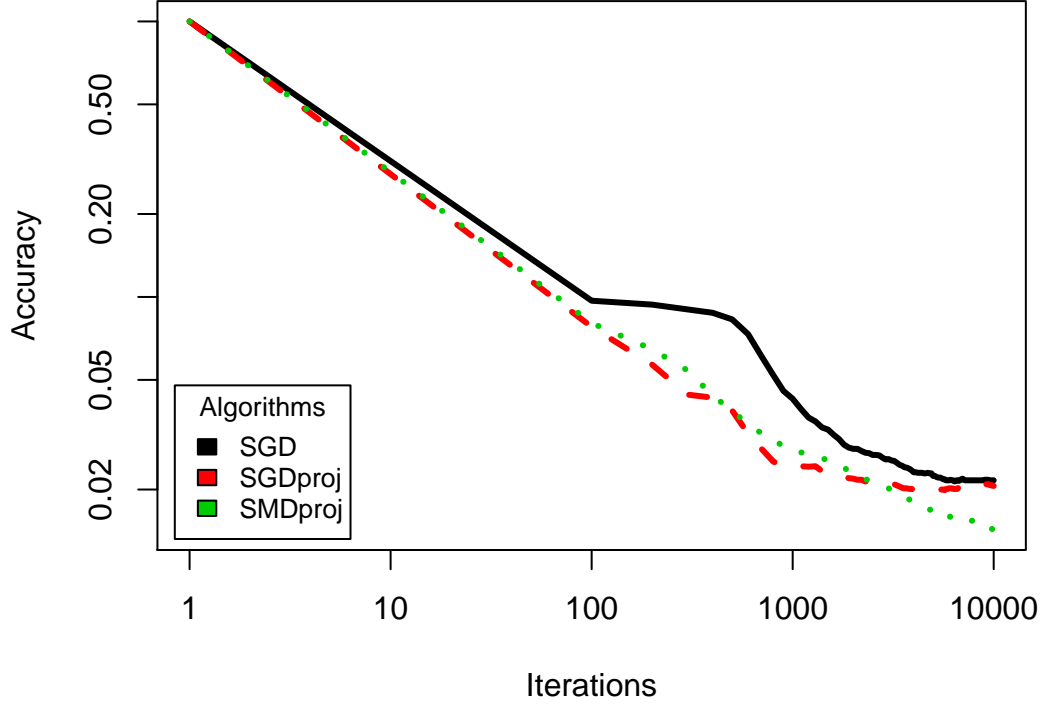
Iteration: Update

$$\begin{aligned} \eta_t &= 1/\sqrt{t}, \\ y_{t+1} &= y_t - \eta_t \nabla \ell_{a_{I_t}, b_{I_t}}(x_t), \\ x_{t+1} &= \Pi_{B_1(z)}(y_{t+1}). \end{aligned}$$

Return: $\bar{x}_{T+1} = \frac{1}{T+1} \sum_{t=1}^{T+1} x_t$

One implements the stochastic version of OMD on MNIST. Note that the regularization parameter is not required since the OMD presented here solves any convex problem and not only the strongly convex ones. Despite slower theoretical rates of convergences the accuracies are very similar to regularized SGD. The better convergences at the end of the experiments are due to the raise of the regularization that deviates regularized SGD from its objective.

SVM on Test Set from MNIST



3.3.2 Randomized strategies, expert advice

Recall the randomized strategy from the Rock, Paper and Scissor game. More generally, consider the setting of d experts with losses $\ell_{t,i}$, $1 \leq i \leq d$.

Definition 12. *The Expert Advice is the assignment of confidants weights $x_{t,i}$ to each experts $1 \leq i \leq d$ in order to get the best randomized strategy that picks randomly an expert I_t with probability x_{t,I_t} . The aim is to bound the averaged regret*

$$\mathbb{E}[\text{Regret}_T(\ell)] = \mathbb{E} \left[\sum_{t=1}^T \ell_{t,I_t} - \min_{1 \leq i \leq d} \sum_{t=1}^T \ell_{t,i} \right].$$

We notice that

$$\mathbb{E}[\text{Regret}_T(\ell)] = \sum_{t=1}^T \mathbb{E}_{x_t}[\ell_{t,I_t}] - \min_{1 \leq i \leq d} \sum_{t=1}^T \ell_{t,i}.$$

Denoting the linear loss function $f_t(x_t) = \mathbb{E}_{x_t}[\ell_{t,I_t}] = \sum_{i=1}^d x_{t,i} \ell_{t,i} = x_t^T \ell_t$ it is an OCO problem on the linear loss over the simplex Λ .

Exercise 10. *Check that $\min_{x \in \Lambda} x^T \sum_{t=1}^T \ell_t = \min_{1 \leq i \leq d} \sum_{t=1}^T \ell_{t,i}$.*

Let $R(x) = x^T \log(x) = \sum_{i=1}^d x_i \log(x_i)$ be the negative entropy function. We consider it as a regularization function over Λ since

$$\nabla R(x) = 1 + \log(x), \quad \nabla^2 R(x) = \text{Diag}(1/x^2) \succeq I_d$$

even if it is not well defined on the boundary of the simplex when $x_i = 0$ for some $1 \leq i \leq d$. Note the the dual space $\{\nabla R(x), x \in \text{Dom}(\nabla R)\} = \mathbb{R}^d$. In order to express OMD (lazy version), we obtain an expression for

$$\nabla R^*(y) = \arg \max_{x \in \Lambda} \{y^T x - R(x)\}.$$

We compute the Lagrangian function

$$\mathcal{L}(x, \theta, \zeta) = y - \nabla R(x) + \theta \left(\sum_{i=1}^d x_i - 1 \right) - \sum_{i=1}^d \zeta_i x_i$$

with parameters $x \in \mathbb{R}^d$, $\theta \in \mathbb{R}$ and $\zeta \in \mathbb{R}_+^d$. We deduce its gradient

$$\nabla \mathcal{L}(x, \theta, \zeta) = \begin{pmatrix} y - \log(x) + (\theta - 1) \mathbb{1} - \zeta \\ \sum_{i=1}^d x_i - 1 \\ -x \end{pmatrix}, \quad \mathbb{1} = (1, \dots, 1)^T.$$

From KKT, we get the conditions

$$\begin{cases} x & = \exp(y + (\theta - 1) \mathbb{1} - \zeta) \\ \sum_{i=1}^d x_i & = 1 \\ x_i \zeta_i & = 0, \quad 1 \leq i \leq d. \end{cases}$$

From the first condition we see that $x_i > 0$ so that ζ must be null. We get

$$x = \frac{\exp(y)}{\sum_{i=1}^d \exp(y_i)}.$$

We obtain the randomized strategy of expert advice called Exponentiated Weighting Algorithm.

Algorithm 11: EWA, Littlestone and Warmuth (1994)

Parameters: step-size $\eta > 0$.

Initialization: Initial prediction $x_1 = (1/d) \mathbb{1}$ and $\theta_1 = 0$.

For each recursion $t \geq 1$:

Sample an expert: $I_t \sim x_t$

Predict as the I_t -th expert

Incur the loss: ℓ_{t, I_t}

Observe: $\ell_t \in \mathbb{R}^d$

Recursion: Update

$$\begin{aligned} \theta_{t+1} &= \theta_t - \eta \ell_t, \\ x_{t+1} &= \frac{\exp(\theta_{t+1})}{\sum_{i=1}^d \exp(\theta_{t+1, i})}. \end{aligned}$$

Exercise 11. Show that EWA coincides with the agile version of OMD with negative entropic regularizer.

Remark 4. We have complete information since we observe all the losses $\ell_t \in \mathbb{R}^d$ despite we pick only one expert $I_t \sim x_t$.

We note that $\nabla f_t(x_t) = \ell_t$ because $f_t(x_t) = x_t^T \ell_t$. Since $\nabla R^*(\theta) = \exp(\theta) / (\sum_{i=1}^d \exp(\theta_i)) = x$ we get $\nabla^2 R^*(y) = \text{Diag}(x) - xx^T$ so that

$$\|\nabla f_t(x_t)\|_t^{*2} = \|\ell_t\|_{\nabla^2 R^*(\bar{\theta})}^2$$

for some $\tilde{\theta}$ in the segment $[\theta_t, \theta_{t+1}]$ so that denoting $\tilde{x} = \exp(\tilde{\theta}) / (\sum_{i=1}^d \exp(\tilde{\theta}_i))$ the corresponding weights, we obtain

$$\|\nabla f_t(x_t)\|_t^{*2} = \left(\sum_{i=1}^d \tilde{x}_i \ell_{t,i}^2 - \left(\sum_{i=1}^d \tilde{x}_i \ell_{t,i} \right)^2 \right) \leq G_\infty^2$$

where $|\ell_{t,i}| \leq G_\infty$ for all $t \geq 1$, $1 \leq i \leq d$. We obtain, denoting $f(x) = x^T \ell$ that

$$\mathbb{E}[\text{Regret}_T(\ell)] = \text{Regret}_T(f) \leq \frac{\eta T G_\infty^2}{2} + \frac{\log d}{\eta} \leq G_\infty \sqrt{2T \log d},$$

choosing $\eta = G_\infty^{-1} \sqrt{2 \log d / T}$.

Remark 5. *The dependence on the dimension in $\sqrt{\log d}$ is optimal and is due to the use of the duality and the ℓ^1 -ball. Indeed $G_\infty = \max_{1 \leq i \leq d} |\ell_{t,i}| = \sup_{x \in \Lambda} \|\nabla f_t(x)\|_\infty = G_{R^*}$ and*

$$\log d = \sup_{x, x' \in \Lambda} R(x) - R(x') = D_R^2$$

the "diameter" of Λ for the regularization function. Thus the regret bound is of order $G_{R^} D_R \sqrt{T}$. Here the dependence in the dimension of $G_{R^*} D_R$ is optimal in $\sqrt{\log d}$. Such a dependence in the dimension is not achievable when constraining to the ℓ^2 -ball whatever is the choice of R . There the dependence in $GD \approx \sqrt{d}$ in the regret analysis of OGD is optimal without additional constrained.*

Example 4. *Back to Rock, Paper and Scissor we get an averaged regret bound which is sub-linear in the complete adversarial setting*

$$\sum_{t=1}^T \max_{y \in \Lambda} y^T A x_t - \min_{1 \leq i \leq 3} \sum_{t=1}^T \max_{y \in \Lambda} y^T A_i \leq \sqrt{2T \log 3},$$

where A_i represent the i -th column of the cost matrix. Since any deterministic strategy incur a loss of 1 at each round in the complete adversarial setting, we get

$$\sum_{t=1}^T \max_{y \in \Lambda} y^T A x_t \leq T + \sqrt{2T \log 3},$$

It is a useless bound and shows the limit of the regret analysis which is relative to a fixed strategy that can be bad in a complete adversarial setting.

A powerful consequence of this analysis is the combination of EWA with the gradient trick for the OCO problem on $\mathcal{K} = \Lambda$. For any algorithm we have the gradient trick

$$\text{Regret}_T(f) \leq \sum_{t=1}^T \nabla f_t(x_t) (x_t - x^*) \leq \sum_{t=1}^T (\nabla f_t(x_t)^T x_t - \nabla f_t(x_t)^T x^*).$$

One interprets the upper bound as the regret of a randomized strategy of d experts with linearized losses $\ell_t = \nabla f_t(x_t)$. We obtain

Algorithm 12: Hedge, Littlestone and Warmuth (1994)

Parameters: step-size $\eta > 0$.

Initialization: Initial prediction $x = (1/d)\mathbb{1}$ and $\theta_1 = 0$.

For each recursion $t \geq 1$:

Predict: x_t

Incur: $f_t(x_t)$

Observe: $\nabla f_t(x_t) \in \mathbb{R}^d$

Recursion: Update

$$\begin{aligned}\theta_{t+1} &= \theta_t - \eta \nabla f_t(x_t), \\ x_{t+1} &= \frac{\exp(\theta_{t+1})}{\sum_{i=1}^d \exp(\theta_{t+1,i})}.\end{aligned}$$

We immediately get the optimal regret bound

$$\text{Regret}_T(f) = \sum_{t=1}^T f_t(x_t) - \min_{x \in \Lambda} \sum_{t=1}^T f_t(x) \leq \mathbb{E}[\text{Regret}_T(\nabla f)] \leq G_\infty \sqrt{2T \log d}.$$

Note that thanks to the gradient trick the regret bound is now relative to the best fixed strategy of the simplex.

Example 5. Back to Rock, Paper and Scissor Hedge get an averaged regret bound which is sub-linear in the complete adversarial setting

$$\sum_{t=1}^T \max_{y \in \Lambda} y^T A x_t - \min_{x \in \Lambda} \sum_{t=1}^T \max_{y \in \Lambda} y^T A x \leq \sqrt{2T \log 3},$$

where A_i represent the i -th column of the cost matrix. Since the randomized strategy $(1/3, 1/3, 1/3)$ incurs the loss 0 at each round in the complete adversarial setting, we get

$$\sum_{t=1}^T \max_{y \in \Lambda} y^T A x_t \leq \sqrt{2T \log 3}.$$

It is a very useful bound to prove the von Neumann minimax theorem in zero-sum games.

Exercise 12. In Rock, Paper and Scissor show that the algorithm Hedge coincides with the optimal randomized strategy $(1/3, 1/3, 1/3)$. Compare with the algorithm EWA.

An extension to the OCO over $\mathcal{K} = B_1(z)$ the ℓ^1 -ball of radius $z > 0$ is achieved using $2d$ -experts from the following representation

Lemma 1. Every $x \in B_1(z)$ satisfies $x_i = z(w_i - w_{i+d})$, $1 \leq i \leq d$, where $w \in \Lambda_{2d}$.

Proof. Introduce d parameters $\lambda_i \geq 0$ and define

$$w_i = x_{i+}/z + \lambda_i, \quad w_{i+d} = x_{i-}/z + \lambda_i, \quad 1 \leq i \leq d,$$

where $x_{i\pm}/z = \max(\pm x_i, 0)$. Then $x_i = z(w_i - w_{i+d})$ and $w \in \Lambda_{2d}$ if and only if

$$\|w\|_1 = \sum_{i=1}^{2d} w_i = \|x\|_1/z + \|\lambda\|_1 = 1.$$

There exists such $\lambda \geq 0$ because $\|x\|_1/z \leq 1$ (even infinitely many when $\|x\|_1 < z$). \square

Indeed the gradient trick still holds and denoting $x_i = z(w_i - w_{i+d})$ the linearized loss

$$\nabla f_t(x_t)^T x_t = \sum_{i=1}^d \nabla f_t(x_t)_i x_{t,i} = z \left(\sum_{i=1}^d \nabla f_t(x_t)_i w_{t,i} - \sum_{i=1}^d \nabla f_t(x_t)_i w_{t,i+d} \right) = z \pm \widehat{\nabla f_t}(x_t)^T w_t$$

where $w_t \in \Lambda_{2d}$ and

$$\pm \widehat{\nabla f_t}(x_t) = (\nabla f_t(x_t)_1, \dots, \nabla f_t(x_t)_d, -\nabla f_t(x_t)_1, \dots, -\nabla f_t(x_t)_d) \in \mathbb{R}^{2d}.$$

Thus we obtain a reduction of the OCO problem on $B_1(z)$ to the OCO problem on $w \in \Lambda_{2d}$ and we get

Algorithm 13: Exponentiated Gradient +/-, Kivinen and Warmuth (1997)

Parameters: step-size $\eta > 0$, radius $z > 0$.

Initialization: Initial prediction $x = 0$ weights $w = 1/(2d) \mathbb{1}$ and $\theta_1 = 0 \in \mathbb{R}^{2d}$.

For each recursion $t \geq 1$:

Predict: x_t

Incur: $f_t(x_t)$

Observe: $\nabla f_t(x_t) \in \mathbb{R}^d$

Recursion: Update

$$\begin{aligned} \theta_{t+1,i} &= \theta_{t,i} - \eta \nabla f_t(x_t)_i, & 1 \leq i \leq d, \\ \theta_{t+1,i} &= \theta_{t,i} + \eta \nabla f_t(x_t)_i, & d+1 \leq i \leq 2d, \\ w_{t+1} &= \frac{\exp(\theta_{t+1})}{\sum_{i=1}^{2d} \exp(\theta_{t+1,i})}, \\ x_{t+1,i} &= z(w_{t+1,i} - w_{t+1,i+d}), & 1 \leq i \leq d. \end{aligned}$$

We immediately get the optimal regret bound choosing $\eta = (zG_\infty)^{-1} \sqrt{2 \log d/T}$

$$\text{Regret}_T \leq G_\infty z \sqrt{2T \log(2d)}.$$

One implements the stochastic version of EG +/- on MNIST and improved the performances of SMD (the radius of the ℓ^1 -ball is still $z = 100$).

Algorithm 14: SEG+/- for linear SVM

Parameters: Epoch T , radius $z > 0$.

Initialization: Initial point $x_1 = 0$, weights $w = 1/(2d) \mathbb{1}$ and $\theta_1 = 0 \in \mathbb{R}^{2d}$.

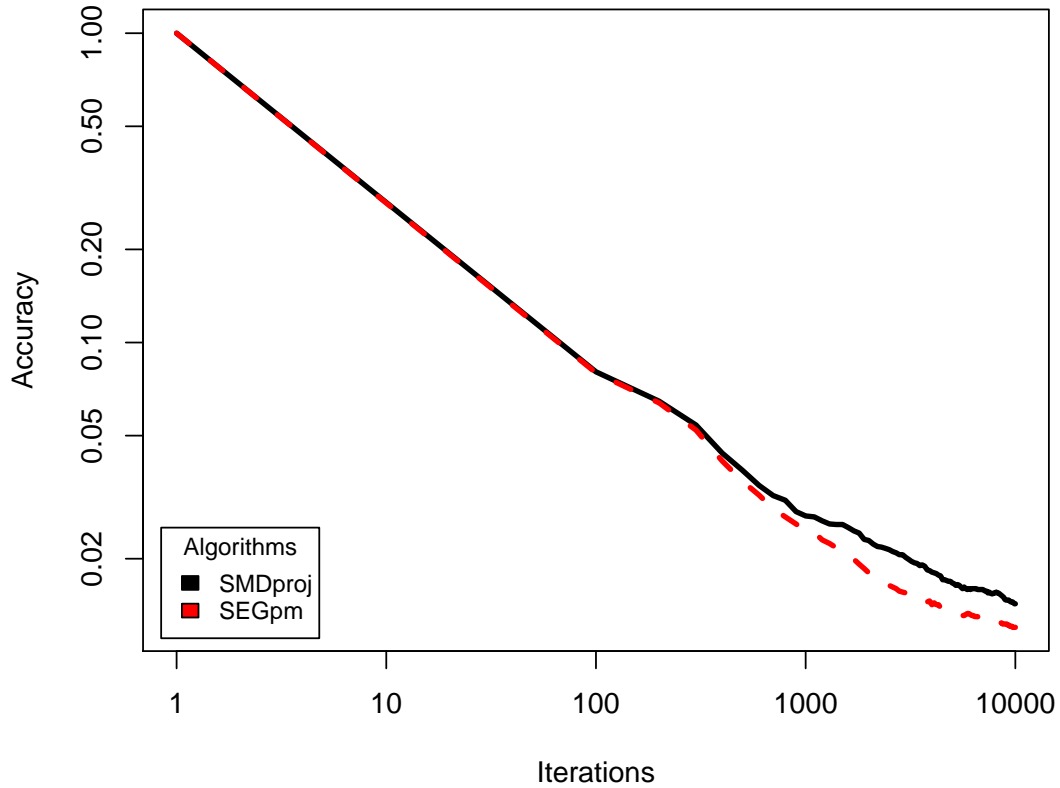
Sample uniformly iid: $(I_t)_{1 \leq t \leq T}$ from $\{1 \leq i \leq n\}$

For each iteration $t = 1, \dots, T$:

Iteration: Update

$$\begin{aligned} \eta_t &= \sqrt{1/t} \\ \theta_{t+1} &= \theta_t - \eta_t \pm \nabla \ell_{a_{I_t}, b_{I_t}}(x_t), \\ w_{t+1} &= \frac{\exp(\theta_{t+1})}{\sum_{i=1}^{2d} \exp(\theta_{t+1,i})}, \\ x_{t+1,i} &= z(w_{t+1,i} - w_{t+1,i+d}), \quad 1 \leq i \leq d. \end{aligned}$$

Return: $\bar{x}_{T+1} = \frac{1}{T+1} \sum_{t=1}^{T+1} x_t$

SVM on Test Set from MNIST

3.3.3 AdaGrad

We recall the regret bound for the general OMD or RFTL

$$\text{Regret}_T(u) \leq \frac{\eta}{2} \sum_{t=1}^T \|\nabla f_t(x_t)\|_t^{*2} + \frac{R(u) - R(x_1)}{\eta}$$

which is equal to, if η is optimized by an oracle,

$$\text{Regret}_T(u) \leq \sqrt{2 \sum_{t=1}^T \|\nabla f_t(x_t)\|_t^{*2} (R(u) - R(x_1))}.$$

As we saw this bound heavily depends on the choice of the regularization function R . The best choice of R heavily depends on the properties of the gradients $\nabla f_t(x_t)$ of the losses of the algorithm itself. AdaGrad will learn how to choose the best regularization function.

We restrict R to the class of weighted quadratic regularization functions $R \in \mathcal{H}$ satisfying

$$\forall x \in \mathcal{K}, \nabla^2 R(x) = D = \text{Diag}(s), \quad s \in (0, \infty)^d, \|s\|_1 \leq 1, .$$

Remark 6.

- $R(x) = \frac{1}{2d} \|x\|^2$ such that $\nabla^2 R(x) = \frac{1}{d} I_d$ and $R \in \mathcal{H}$,
- $R(x) = x^T \log(x)$ such that $\nabla^2 R(x) = \text{Diag}(1/x)$ is not in \mathcal{H} .

We first determine what could be the best possible regret bound. We compute the second derivative of the convex conjugate $\nabla^2 R^*(x^*) = D^{-1}$ for $D \succ 0$ and

$$\begin{aligned} \min_{R \in \mathcal{H}} \sum_{t=1}^T \|\nabla f_t(x_t)\|_t^{*2} &= \min_{D = \text{Diag}(s)} \sum_{t=1}^T \|\nabla f_t(x_t)\|_{D^{-1}}^2 \\ &= \min_{s \in \mathbb{R}_+^d, \|s\|_1 \leq 1} \sum_{t=1}^T \sum_{i=1}^d \nabla f_t(x_t)_i^2 s_i^{-1} \end{aligned}$$

Applying Cauchy-Schwartz, we have

$$\sum_{i=1}^d \left(\sqrt{\sum_{t=1}^T \nabla f_t(x_t)_i^2 / s_i} \right)^2 \sum_{i=1}^d \sqrt{s_i} \geq \left(\sum_{i=1}^d \sqrt{\sum_{t=1}^T \nabla f_t(x_t)_i^2} \right)^2$$

so that

$$\min_{s \in \mathbb{R}_+^d, \|s\|_1 \leq 1} \sum_{t=1}^T \sum_{i=1}^d \nabla f_t(x_t)_i^2 s_i^{-1} \geq \left(\sum_{i=1}^d \sqrt{\sum_{t=1}^T \nabla f_t(x_t)_i^2} \right)^2.$$

Note that this minimizer is achieved by $\|s^*\|_1 = 1$ for

$$s_i^* = \frac{\sqrt{\sum_{t=1}^T \nabla f_t(x_t)_i^2}}{\sum_{i=1}^d \sqrt{\sum_{t=1}^T \nabla f_t(x_t)_i^2}}.$$

Thus the best possible regret in this class of regularization function is

$$\text{Regret}_T(u) \leq \sum_{i=1}^d \sqrt{2(R(u) - R(x_1)) \sum_{t=1}^T \nabla f_t(x_t)_i^2}.$$

However such learning rate is not known before step T . AdaGrad solves this problem by considering a **multiple adaptive learning rates** approach; each coordinate will have its own gradient step close to the optimal s_i .

Algorithm 15: AdaGrad (diagonal version), Duchi et al. (2011)

Parameters: step-size $\eta > 0$.

Initialization: Initial prediction $y_1 = x_1 \in \mathcal{K}$, initial multiple learning rates $S_0 = 0$ (or $= \delta \mathbb{I}$ small).

Predict: x_t

Incur the average loss: $f_t(x_t)$

Observe: $\nabla f_t(x_t) \in \mathbb{R}^d$

Recursion: Update

$$\begin{aligned} S_t &= S_{t-1} + \nabla f_t(x_t)^2 \\ D_t &= \text{Diag}(\sqrt{S_t}) \\ y_{t+1} &= x_t - \eta D_t^{-1} \nabla f_t(x_t), \\ x_{t+1} &= \arg \min_{x \in \mathcal{K}} \|x - y_{t+1}\|_{D_t}^2, \quad 1 \leq i \leq d. \end{aligned}$$

We notice that AdaGrad is an agile OMD algorithm with adaptive regularization functions $R_t(x) = \frac{1}{2} \|x - x_1\|_{D_t}^2$ since then

$$\nabla R_t(y_{t+1}) = \nabla R_t(x_t) - \eta \nabla f_t(x_t), \quad B_{R_t}(x|y) = \frac{1}{2} \|x - y\|_{D_t}^2.$$

Theorem 10. For $\eta = D_\infty / \sqrt{2}$ with $D_\infty = \max_{x,y \in \mathcal{K}} \|x - y\|_\infty$ AdaGrad get the regret bound

$$\text{Regret}_T \leq D_\infty \sum_{i=1}^d \sqrt{2 \sum_{t=1}^T \nabla f_t(x_t)_i^2}.$$

Since the regularization functions are depending on t , one has to adapt the simple analysis of OGD above.

Proof. We start from the recursive relation $y_{t+1} - u = y_t - u - \eta D_t^{-1} \nabla f_t(x_t)$ that we rewrite as $D_t(y_{t+1} - u) = D_t(y_t - u) - \eta \nabla f_t(x_t)$ so that multiplying both relations we get

$$\begin{aligned} \|y_{t+1} - u\|_{D_t}^2 &= (y_{t+1} - u)^T D_t (y_{t+1} - u) \\ &= \|x_t - u\|_{D_t}^2 - 2\eta \nabla f_t(x_t)^T (x_t - u) + \eta^2 \|\nabla f_t(x_t)\|_{D_t^{-1}}^2. \end{aligned}$$

By the pythagorean Theorem we also have $\|x_{t+1} - u\|_{D_t}^2 \leq \|y_{t+1} - u\|_{D_t}^2$ so that

$$2\eta \nabla f_t(x_t)^T (x_t - u) \leq \|x_t - u\|_{D_t}^2 - \|x_{t+1} - u\|_{D_t}^2 + \eta^2 \|\nabla f_t(x_t)\|_{D_t^{-1}}^2.$$

Then we get

$$\begin{aligned} 2 \sum_{t=1}^T \nabla f_t(x_t)^T (x_t - u) &\leq \frac{1}{\eta} \sum_{t=1}^T \left(\|x_t - u\|_{D_t}^2 - \|x_{t+1} - u\|_{D_t}^2 + \eta^2 \|\nabla f_t(x_t)\|_{D_t^{-1}}^2 \right) \\ &\leq \frac{1}{\eta} \sum_{t=1}^T \left(\|x_t - u\|_{D_t}^2 - \|x_t - u\|_{D_{t-1}}^2 \right) \\ &\quad + \eta \sum_{t=1}^T \|\nabla f_t(x_t)\|_{D_t^{-1}}^2, \end{aligned}$$

with the convention $D_0 = \text{Diag}(S_0)$.

For the first term we use the telescoping sum

$$\sum_{t=1}^T (x_t - u)^T (D_t - D_{t-1})(x_t - u) \leq D_\infty^2 \sum_{i=1}^d \sum_{t=1}^T (\sqrt{S_t} - \sqrt{S_{t-1}})_i \leq D_\infty^2 \sum_{i=1}^d \sqrt{S_{T,i}}.$$

For the last term, we get

$$\begin{aligned} \sum_{t=1}^T \|\nabla f_t(x_t)\|_{D_t^{-1}}^2 &= \sum_{t=1}^T \sum_{i=1}^d (\nabla f_t(x_t)^2 / \sqrt{S_t})_i \\ &\leq \sum_{i=1}^d \sum_{t=1}^T ((S_t - S_{t-1}) / \sqrt{S_t})_i. \end{aligned}$$

By a comparison with an integral, we get

$$\sum_{t=1}^T (S_{t,i} - S_{t-1,i}) / \sqrt{S_{t,i}} \leq \int_0^{S_{T,i}} \frac{dx}{\sqrt{x}} = 2\sqrt{S_{T,i}}.$$

Finally we obtain the regret bound

$$\frac{1}{2} \left(\frac{D_\infty^2}{\eta} + 2\eta \right) \sum_{i=1}^d \sqrt{S_{T,i}}$$

which yields the desired result as $\eta = D_\infty / \sqrt{2}$. \square

Remark 7. Denoting $D_i = \max_{x,y \in \mathcal{K}} \|(x - y)_i\|$ we immediately improve the regret bound to

$$\text{Regret}_T \leq \sum_{i=1}^d D_i \sqrt{2 \sum_{t=1}^T \nabla f_t(x_t)_i^2} \leq \sqrt{\sum_{i=1}^d D_i^2} \sqrt{2 \sum_{t=1}^T \nabla f_t(x_t)_i^2}$$

by Cauchy-Schwartz inequality. Since for hyper-rectangles \mathcal{K} we have $\sum_{i=1}^d D_i^2 = D^2 = \max_{x,y \in \mathcal{K}} \|x - y\|^2$, if the gradients were the same for OGD and Adagrad then the best possible regret for OGD is always larger than the one for Adagrad.

In order to implement AdaGrad to the linear SVM one has to explicit the projection step.

Algorithm 16: Adagrad for linear SVM

Parameters: Epoch T , radius $z > 0$.

Initialization: Initial point $x_1 = y_1 = 0$ and $S_0 = 0$ (or $= \delta \mathbb{I}$ small).

Sample uniformly iid: $(I_t)_{1 \leq t \leq T}$ from $\{1 \leq i \leq n\}$

For each iteration $t = 1, \dots, T$:

Iteration: Update

$$\begin{aligned} S_t &= S_{t-1} + \nabla \ell_{a_{I_t}, b_{I_t}}(x_t)^2 \\ D_t &= \text{Diag}(\sqrt{S_t}) \\ y_{t+1} &= x_t - D_t^{-1} \nabla \ell_{a_{I_t}, b_{I_t}}(x_t), \\ x_{t+1} &= \arg \min_{x \in B_1(z)} \|x - y_{t+1}\|_{D_t}^2, \quad 1 \leq i \leq d. \end{aligned}$$

Return: $\bar{x}_{T+1} = \frac{1}{T+1} \sum_{t=1}^{T+1} x_t$

For that one has to adapt the Euclidian projection on the simplex to weighted norms. One has to solve the CO

$$\arg \min_{w \in B_1(z)} \sum_{i=1}^d \|w - x\|_D^2.$$

We have the Lagrangian function

$$\mathcal{L}(x, \theta, \zeta) = \frac{1}{2}(w - x)^T D(w - x) + \theta \left(\sum_{i=1}^d w_i - 1 \right) - \sum_{i=1}^d \zeta_i w_i$$

with parameters $w \in \mathbb{R}^d$, $\theta \in \mathbb{R}$ and $\zeta \in \mathbb{R}_+^d$. We compute its gradient

$$\nabla \mathcal{L}(w, \theta, \zeta) = \begin{pmatrix} D(w - x) + \theta \mathbb{I} - \zeta \\ \sum_{i=1}^d w_i - 1 \\ -w \end{pmatrix}, \quad \mathbb{I} = (1, \dots, 1)^T.$$

Thus KKT provides

$$\begin{cases} w^* = x - D^{-1}(\theta^* \mathbb{I} + \zeta^*), \\ \sum_{i=1}^d w_i^* = 1 \\ w_i^* = 0 \text{ or } w_i^* > 0 \text{ and } \zeta_i^* = 0. \end{cases}$$

To sum up we obtain the weighted soft-thresholding

$$w_i^* = \max(x - D^{-1}\theta^* \mathbb{I}, 0) = D^{-1} \text{SoftThreshold}(Dx, \theta^*).$$

Thus denoting $\|w^*\|_0 = d_0$ we get the relation

$$1 = \sum_{j=1}^{d_0} w_{(j)}^* = \sum_{j=1}^{d_0} D^{-1} \text{SoftThreshold}(Dx, \theta^*) = \sum_{j=1}^{d_0} x_{(j)} - \sum_{j=1}^{d_0} D_{(j)}^{-1} \theta^*$$

where $D_{(j)}$ is the diagonal element of D with the same ordering so that necessarily

$$\theta^* = \frac{1}{\sum_{j=1}^{d_0} D_{(j)}^{-1}} \left(\sum_{j=1}^{d_0} x_{(j)} - 1 \right).$$

We obtain

Algorithm 17: Projection on the simplex with weighted norm $\|\cdot\|_D$

Input: $x \in \mathbb{R}^d$ and D diagonal.

If $x \in \Lambda$

Then Return x .

Else

Sort $(Dx)_{(1)} \geq \dots \geq (Dx)_{(d)}$

Find $d_0 = \max \left\{ 1 \leq i \leq d; (Dx)_{(i)} - \frac{1}{\sum_{j=1}^i D_{(j)}^{-1}} (\sum_{j=1}^i x_{(j)} - 1) \right\}$

Define $\theta^* = \frac{1}{\sum_{j=1}^{d_0} D_{(j)}^{-1}} \left(\sum_{j=1}^{d_0} x_{(j)} - 1 \right)$

Return $w^* = D^{-1} \text{SoftThreshold}(Dx, \theta^*)$.

Recall that for the hinge loss we have, for any instance (a, b) of the training set $(a_t, b_t)_{1 \leq t \leq d}$

$$\nabla \ell_{a,b}(x_t) = \begin{cases} 0 & \text{if } bx^T a > 1, \\ -ba & \text{else.} \end{cases}$$

Assume that the design is sparse, i.e. we have

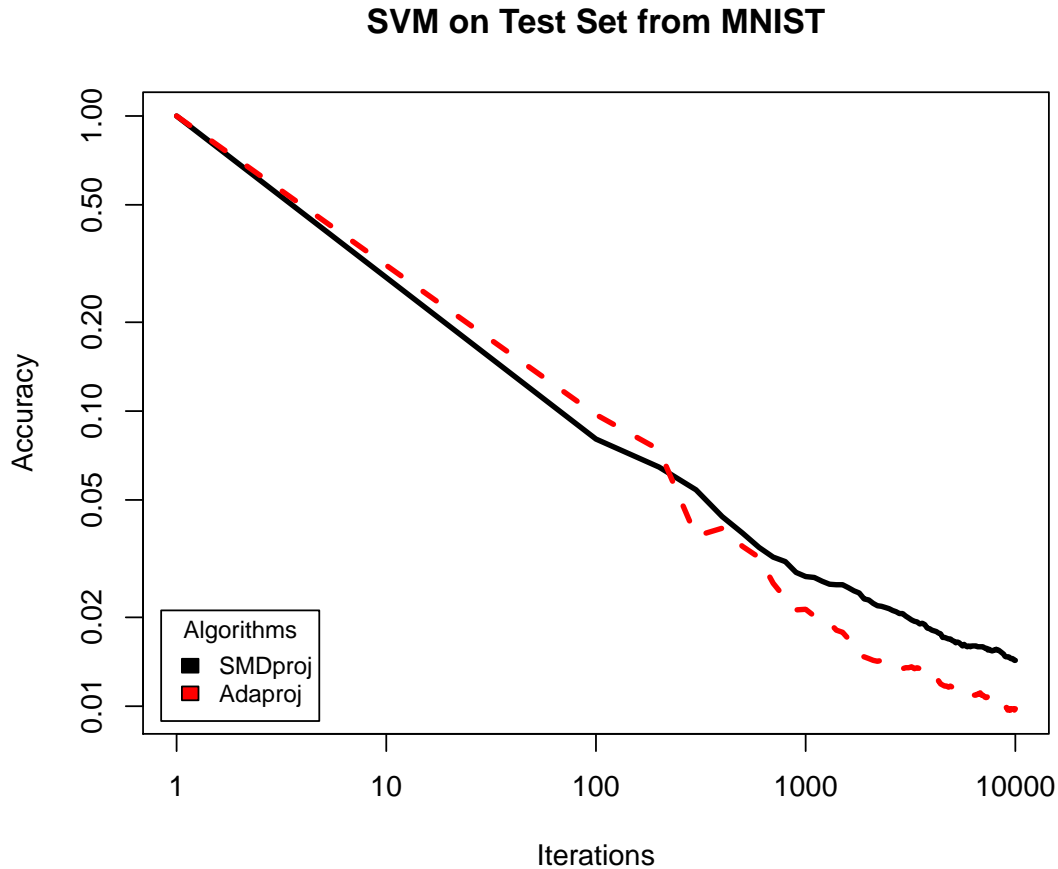
$$\mathbb{P}(a_i \neq 0) = \min\{1, ci^{-\alpha}\}$$

for some $\alpha \in (2, \infty)$ and any $1 \leq i \leq d$. Then the regret of AdaGrad is low in expectation

$$\begin{aligned} \mathbb{E}[\text{Regret}_t] &\leq D_\infty \sum_{i=1}^d \mathbb{E} \left[\sqrt{2 \sum_{t=1}^T \nabla f_t(x_t)_i^2} \right] \\ &\leq D_\infty \sum_{i=1}^d \sqrt{2 \sum_{t=1}^T \mathbb{E} \left[\nabla f_t(x_t)_i^2 \right]} \\ &\leq D_\infty G_\infty \sqrt{2T} \sum_{i=1}^d \sqrt{\mathbb{P}(a_i \neq 0)} \\ &\leq D_\infty G_\infty \sqrt{2cT} \sum_{i=1}^d i^{-\alpha/2} \\ &\leq D_\infty \sqrt{2cT} (1 + \log d), \end{aligned}$$

where $\|\nabla f_t(x_t)\|_\infty \leq G_\infty$ for any $1 \leq t \leq T$. Taking advantage of the sparsity thanks to the adaptivity of AdaGrad one turns a $G \approx \sqrt{d}$ regret bound of OGD into a $\log d$ one.

Implemented on MNIST, AdaGrad clearly takes advantage of the sparsity in the pixels of the the handwritten digits in a better way than the projection of the ℓ^1 -ball. It is due to the fact that AdaGrad learns the sparsity via the gradients whereas the radius of the ℓ^1 -ball (or equivalently the regularization parameter in the dual LASSO problem) is fixed a priori (here arbitrarily to $z = 100$).



3.3.4 BOA

BOA is a multiple learning rate version of EG. The idea is to combine the adaptivity of the simplicity of the gradients as in Adagrad together with the use of the geometry of the convex set $\mathcal{K} = \Lambda$ via the negative entropy regularization function. Note that it is necessary to add a quadratic compensation to the gradient in the exponential weights in order to get theoretical guarantees.

Implemented on MNIST the rate seems to be faster than for Adaproj. It is due to the introduction of the quadratic compensation that can be seen as an estimation of the noise level. Then BOA seeks at achieving a good bias-variance tradeoff in stochastic environment. Note that an even better bias-variance tradeoff will be achieved by Adam thanks to momentum.

Algorithm 18: SBOA+/- for linear SVM, Wintenberger (2017)

Parameters: Epoch T , radius $z > 0$.

Initialization: Initial point $x_1 = 0$, weights $w = 1/(2d) \mathbb{I}$ and $\eta_0 = \mathbb{I} \in \mathbb{R}^{2d}$.

Sample uniformly iid: $(I_t)_{1 \leq t \leq T}$ from $\{1 \leq i \leq n\}$

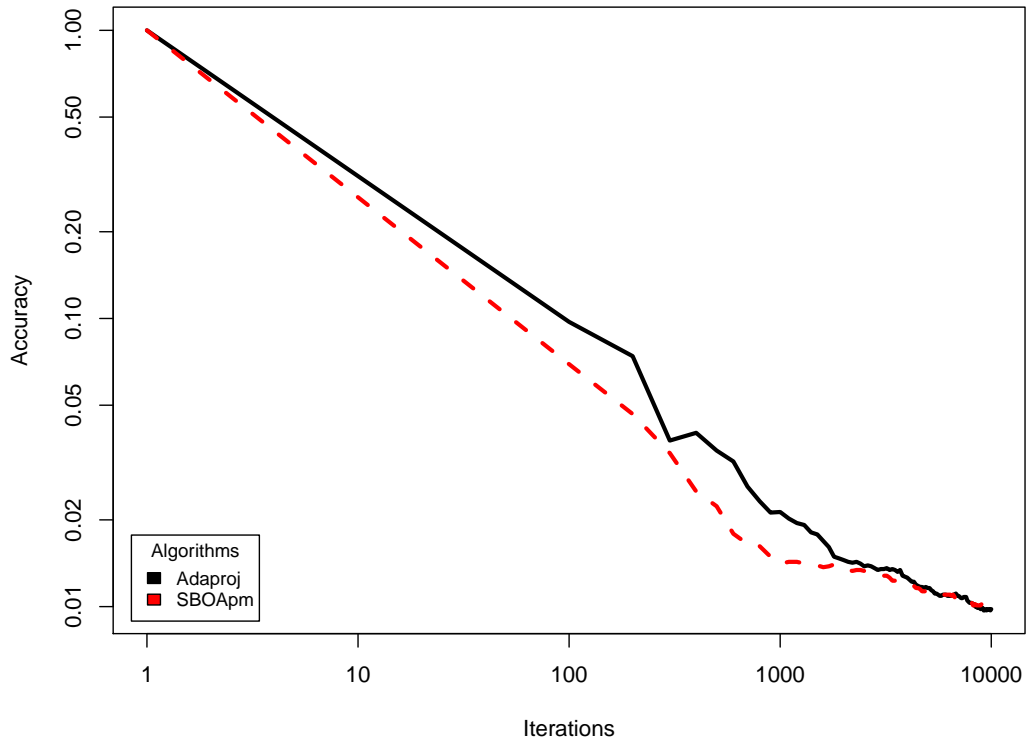
For each iteration $t = 1, \dots, T$:

Iteration: Update

$$\begin{aligned} \bar{\nabla} \ell_t &= w_t^T \pm \nabla \ell_{a_{I_t}, b_{I_t}}(x_t) \\ \theta_{t+1} &= \theta_t - \pm \nabla \ell_{a_{I_t}, b_{I_t}}(x_t) - \eta_{t-1} (\pm \nabla \ell_{a_{I_t}, b_{I_t}}(x_t) - \bar{\nabla} \ell_t)^2, \\ \eta_t &= \sqrt{\eta_{t-1}^2 / (1 + \eta_{t-1}^2 (\pm \nabla \ell_{a_{I_t}, b_{I_t}}(x_t) - \bar{\nabla} \ell_t)^2)} \\ w_{t+1} &= \frac{\eta_t \exp(\eta_t \theta_{t+1})}{\sum_{i=1}^{2d} \eta_{t,i} \exp(\eta_t \theta_{t+1,i})}, \\ x_{t+1,i} &= z(w_{t+1,i} - w_{t+1,i+d}), \quad 1 \leq i \leq d. \end{aligned}$$

Return: $\bar{x}_{T+1} = \frac{1}{T+1} \sum_{t=1}^{T+1} x_t$

SVM on Test Set from MNIST



Part III

Acceleration and exploration

Accelerated OCO algorithms

4.1 Momentum

A variant of AdaGrad with **momentum** is the popular Adam, Kingma and Ba (2014). Momentum has been introduced first by Nesterov (1983) as an acceleration scheme in the CO method. Initially the momentum step was applied to the iterate x_t of the algorithm. It can also accelerate OCO algorithms in practice in the stochastic OCO setting. A natural way of introducing momentum in SGD methods is directly on the successive gradients.

Recall that SGD is based on an unbiased noisy version of the CO problem (f, \mathcal{K}) denoted $\widehat{\nabla}f$. In such a setting $\widehat{\nabla}f(x_t)$ is an unbiased estimator of $\nabla f(x_t)$ since $\widehat{\nabla}f$ is a noisy version of ∇f with mean ∇f . One defines a better estimator of this mean by averaging. However the objective $\nabla f(x_t)$ is evolving through time t .

The momentum estimator m_t is an iterative way of approximating $\nabla f(x_t)$ called Exponential Moving Average, namely

$$m_t = \beta m_{t-1} + (1 - \beta) \widehat{\nabla}f(x_t) \quad \iff \quad m_t = (1 - \beta) \sum_{j=0}^{t-1} \beta^j \widehat{\nabla}f(x_{t-j}).$$

Note that since $(1 - \beta) \sum_{j=0}^{t-1} \beta^j = 1 - \beta^t \neq 1$ one should debiased the momentum m_t by dividing it by $1 - \beta^t$. If successive gradients are pointing to different direction (as they are noisy) then the erratic directions could be averaged and canceled. However if β is too large then past gradients are taken into account in the momentum which might introduce a bias in the estimation of the last gradient $\nabla f(x_t)$.

Indeed the variation of each coordinate $m_{t,i}$

$$v_{t,i} = (1 - \beta) \sum_{j=0}^{t-1} \beta^j (\widehat{\nabla}f(x_{t-j})_i - (1 - \beta) \sum_{j=0}^{t-1} \beta^j \widehat{\nabla}f(x_{t-j})_i)^2$$

satisfies the recursive relation

$$v_{t,i} = (1 - \beta)(v_{t-1,i} + \beta(\widehat{\nabla}f(x_{t-1})_i - m_{t-1,i})^2)$$

so that it is comparable to $(1 - \beta)(\widehat{\nabla}f(x_t)_i - m_{t-1,i})^2$ in a stationary regime when $v_{t,i} \approx v_{t-1,i}$. The variation of noisy $\widehat{\nabla}f$ might then be reduced from a factor $(1 - \beta)$ thanks to momentum.

The choice of β is tricky since the larger β the smaller the variations of m_t but the longer the memory of the momentum. If β is well chosen, a momentum on the stochastic gradients might increase the accuracy of the estimation of the true gradient by reducing the variance and thus accelerating the convergence without deteriorating the stability of the algorithm.

Remark 8. *There exists many other ways of accelerating SGD by improving the estimate of the gradient $\nabla f(x_t)$.*

One way is to use moving averages, i.e. considering

$$m_t = \frac{1}{k} \sum_{i=1}^k \widehat{\nabla} f_t(x_{t-i+1})$$

in the recursion instead of the instantaneous noisy gradient $\widehat{\nabla} f(x_t)$. A large k decreases the variation of the averages of an order k^{-1} when x_{t-i+1} is stable for any $i = 1, \dots, k$ but may increase the bias.

A usual mini-batch scheme is that $x_t = x_{t-k+1}$ for a fixed $k > 0$ and for each $t \in k\mathbb{N} + 1$ we consider

$$m_t = \frac{1}{k} \sum_{i=1}^k \widehat{\nabla} f_{t-i+1}(x_t).$$

Then the variance of the estimator of ∇f is decreased of an order k^{-1} without deteriorating the bias. However it increase the complexity of each gradient step by a factor k . It corresponds to an interpolation between OCO and CO. The question of the optimal choice of k is difficult.

The novelty to Adam is to apply a momentum to the squares of the gradient as well. The motivation comes from the multiple learning rates of Adagrad

$$\frac{1}{\sqrt{t}} \frac{1}{\sqrt{\frac{1}{t} \sum_{k=1}^t \widehat{\nabla} f_t(x_t)_i^2}}, \quad 1 \leq i \leq d,$$

and to interpret it as the multiplication of the learning rate $\frac{1}{\sqrt{t}}$ together with the inverse of an estimator of the noise level

$$\sqrt{\frac{1}{t} \sum_{k=1}^t \widehat{\nabla} f_t(x_t)_i^2} \approx \sqrt{\mathbb{E}[\widehat{\nabla} f_t(x_t)_i^2]}.$$

In this interpretation the noise level of the approximation $\widehat{\nabla} f$ is measured according to a moment of order 2. The same reasoning as before shows that a momentum might improve the estimation of the second order moments of the noisy gradients for a well chosen

coefficient β .

Algorithm 19: Adam for linear SVM, Kingma and Ba (2014)

Parameters: Epoch T , radius $z > 0$, $\beta_1 = 0.9$ and $\beta_2 = 0.999$.

Initialization: Initial point $x_1 = y_1 = 0$ and $S_0 = 0$ (or $= \delta \mathbb{I}$ small).

Sample uniformly iid: $(I_t)_{1 \leq t \leq T}$ from $\{1 \leq i \leq n\}$

For each iteration $t = 1, \dots, T$:

Iteration: Update

$$\eta_t = 1/\sqrt{t},$$

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) \nabla \ell_{a_{I_t}, b_{I_t}}(x_t)$$

$$S_t = \beta_2 S_{t-1} + (1 - \beta_2) \nabla \ell_{a_{I_t}, b_{I_t}}(x_t)^2$$

$$D_t = \text{Diag}(\sqrt{S_t})$$

$$y_{t+1} = x_t - \eta_t D_t^{-1} m_t,$$

$$x_{t+1} = \arg \min_{x \in B_1(z)} \|x - y_{t+1}\|_{D_t}^2, \quad 1 \leq i \leq d.$$

Return: $\bar{x}_{T+1} = \frac{1}{T+1} \sum_{t=1}^{T+1} x_t$

In practice one chooses $\beta_2 \gg \beta_1$ as the noise level of the gradient directions are thought as more stable than the direction of the gradient. The best bias-variance trade-off is then achieved for $\beta_2 \approx 1$, taking into account a large number of past squared gradients.

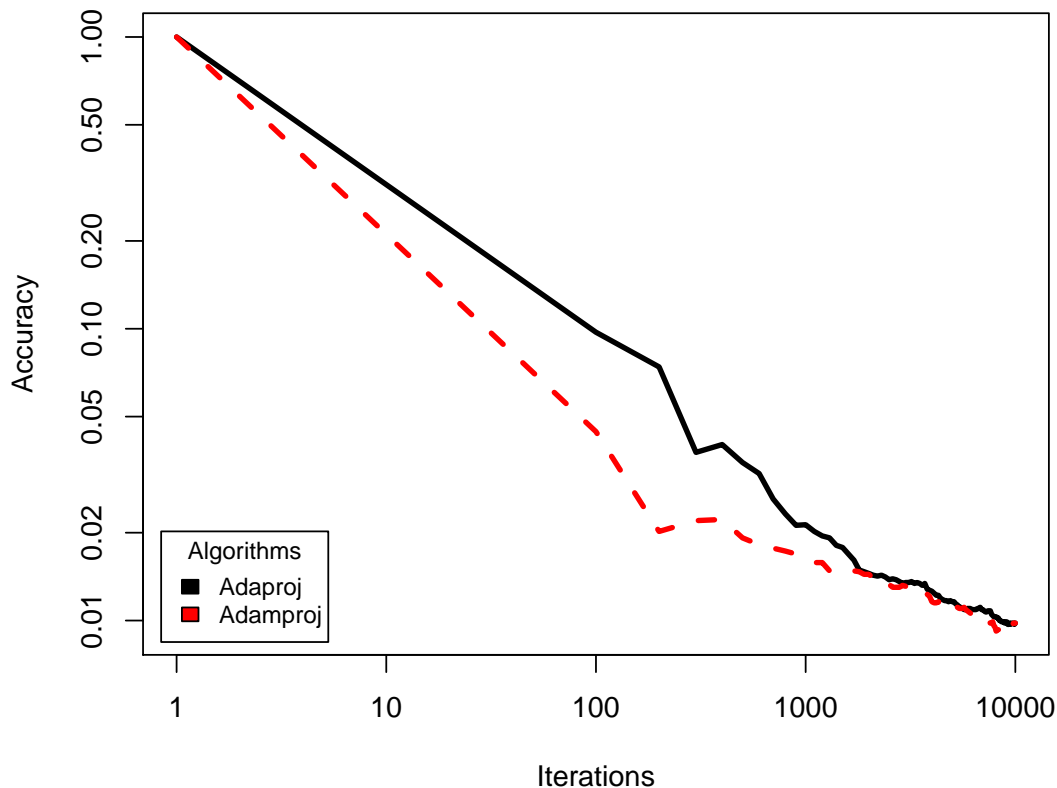
Remark 9. *Acceleration of algorithms might be introduced carefully in order to control the noise stability. Momentum on gradients and on the squared gradients as in Adam may accelerate AdaGrad in optimistic stochastic OCO settings.*

However in an OCO setting with $\mathcal{K} = [-1, 1]$ and $\nabla f_t(x_t) = C > 2$ for $t = 3k - 1$, $k \geq 1$, and $\nabla f_t(x_t) = -1$ otherwise Adam does not converge since $\eta_t D_t^{-1}$ is not a decreasing sequence. Worst Adam converges to 1 whereas $x^ = -1$ for $\beta_2 = (1 + C^2)^{-1}$. Thus β_2 should be taken very large as advised in Kingma and Ba (2014). Worst, for any choice $\sqrt{\beta_2} > \beta_1$, there still exists another OCO setting where the regret of Adam is linear, see Reddi et al. (2019).*

These restrictions remains true for other accelerations such as moving averages discussed above but not for mini-batch.

Adam is very efficient in practice when applied to MNIST dataset. Moreover it works also extremely well in deep learning training, beyond the convex loss function setting (flat high-dimensional problems) where the objective is not necessarily to converge. It explains the success of this algorithm and its variants in deep learning.

SVM on Test Set from MNIST



4.2 Online Newton Step (ONS)

Despite excellent practical acceleration observed in practice in Adam, it is impossible to accelerate OCO algorithms without extra assumptions on the loss function. We have already seen that for α strongly convex loss functions then SGD with learning rates $\eta_t = 1/(\alpha t)$ achieves a logarithmic regret.

4.2.1 Exp-concave functions

In many practical situation the strong convexity assumption is too strong.

Example 6. Consider the linear SVM setting in MNIST with pixels $a \in \mathbb{R}^d$ and label $b \in \{-1, 1\}$ together with the square loss as a relaxation of the 0/1 loss

$$\tilde{\ell}_{a,b}(x) = (b - x^T a)^2.$$

Not that despite the square function $y \rightarrow y^2$ is 2-strongly convex, it is not always the case of $\tilde{\ell}$. Indeed one computes

$$\nabla \tilde{\ell}_{a,b}(x) = 2(b - x^T a)a \quad \text{and} \quad \nabla^2 \tilde{\ell}_{a,b}(x) = 2aa^T,$$

that is convex iff

$$2aa^T \succeq \alpha I_d.$$

We get a contradiction since aa^T is a rank one matrix that cannot be invertible (except for $d = 1$). Thus $\tilde{\ell}_{a,b}$ is not strongly convex.

We have to introduce the notion of exp-concavity

Definition 13. A convex function $f : \mathcal{K} \mapsto \mathbb{R}$ is exp-concave on iff the function $g(x) = \exp(-\mu f(x))$ is concave.

We have the following property

Lemma 2. A twice differentiable function $f : \mathcal{K} \mapsto \mathbb{R}$ is μ -exp-concave iff

$$\nabla^2 f(x) \succeq \mu \nabla f(x) \nabla f(x)^T, \quad x \in \mathcal{K}.$$

Proof. We have $g(x)$ twice differentiable that is concave iff $\nabla^2 g(x) \preceq 0$, $x \in \mathcal{K}$. We compute

$$\begin{aligned} \nabla g(x) &= -\mu \nabla f(x) \exp(-\mu f(x)) \\ \nabla^2 g(x) &= (\mu^2 \nabla f(x) \nabla f(x)^T - \mu \nabla^2 f(x)) \exp(-\mu f(x)) \end{aligned}$$

and $\nabla^2 g(x) \preceq 0$ iff

$$\begin{aligned} \mu^2 \nabla f(x) \nabla f(x)^T - \mu \nabla^2 f(x) &\preceq 0 \\ \mu \nabla f(x) \nabla f(x)^T - \mu \nabla^2 f(x) &\preceq \nabla^2 f(x) \end{aligned}$$

Notice that the rank one matrix $\nabla f(x) \nabla f(x)^T \succeq 0$ by construction so that $\nabla^2 f(x) \succeq 0$ and f is convex. \square

Exercise 13. A α -strongly convex G -Lipschitz function is α/G^2 -exp-concave.

However there are many examples of exp-concave functions that are not strongly convex.

Example 7. In Example 6 we have

$$\nabla \tilde{\ell}_{a,b}(x) \nabla \tilde{\ell}_{a,b}(x)^T = 4(b - x^T a)^2 aa^T \preceq \frac{2}{\mu} aa^T$$

iff $\max_{x \in \mathcal{K}} 2(b - x^T a)^2 \leq \frac{1}{\mu}$. In particular μ is proportional to the amplitude of the square loss and is independent of the dimension of the OCO problem.

We would need the following stronger property of exp-concavity, valid in the usual bounded setting.

Proposition 6. Let $f : \mathcal{K} \mapsto \mathbb{R}$ be μ -exp-concave, D be the diameter of \mathcal{K} and $\max_{x \in \mathcal{K}} \|\nabla f(x)\| \leq G$ for some $G > 0$ as usual. Then

$$f(y) \geq f(x) + \nabla f(x)^T (y - x) + \frac{\gamma}{2} (y - x)^T \nabla f(x) \nabla f(x)^T (y - x), \quad x, y \in \mathcal{K},$$

with $\gamma \leq \frac{1}{2} \min(\frac{1}{4GD}, \mu)$.

Proof. As $2\gamma \leq \mu$ then $h(x) = \exp(-2\gamma f(x))$ is a concave function and

$$h(y) \leq h(x) + \nabla h(x)^T (y - x).$$

We compute $\nabla h(x) = -2\gamma\nabla f(x) \exp(-2\gamma f(x))$ so that by plugging in

$$\begin{aligned} \exp(-2\gamma f(y)) &\leq \exp(-2\gamma f(x)) - \exp(-2\gamma f(x))2\gamma f(x)^T(y-x) \\ &\leq \exp(-2\gamma f(x))(1 - 2\gamma f(x)^T(y-x)). \end{aligned}$$

Thus we get

$$f(y) \geq f(x) - \frac{1}{2\gamma} \log(1 - 2\gamma\nabla f(x)^T(y-x)).$$

Using the boundedness $|2\gamma\nabla f(x)^T(y-x)| \leq 2\gamma GD \leq 1/4$ and that

$$-\log(1-z) \geq z + \frac{1}{4}z^2 \quad |z| \leq 1/4,$$

we obtain

$$f(y) \geq f(x) + \nabla f(x)^T(y-x) + \frac{1}{8\gamma}(2\gamma\nabla f(x)^T(y-x))^2$$

and the desired result follows. \square

4.2.2 Online Newton Step (ONS) algorithm

The ONS is an OCO adaptation of the Newton-Raphson step from CO problems

$$x_{t+1} = x_t - \eta H_t^{-1} \nabla f(x_t)$$

where $H_t = \nabla^2 f(x_t)$ and $\eta > 0$.

In the OCO setting, one can advantageously replace H_t by an approximation as a function of the gradients only, namely $\frac{1}{t} \sum_{k=1}^t \nabla f_k(x_k) \nabla f_k(x_k)^T$, under weaker assumption, namely exp-concavity. We obtain the ONS algorithm

Algorithm 20: Online Newton Step, Hazan and Kale (2011)

Initialization: $\gamma > 0$ and $\epsilon > 0$.

Initialization: Initial prediction $x_1 \in \mathcal{K}$ and $A_0 = \epsilon I_d$.

Predict: x_t

Incur: $f_t(x_t)$

Observe: $\nabla f_t(x_t) \in \mathbb{R}^d$

Recursion: Update

$$\begin{aligned} A_t &= A_{t-1} + \nabla f_t(x_t) \nabla f_t(x_t)^T \\ y_{t+1} &= x_t - \frac{1}{\gamma} A_t^{-1} \nabla f_t(x_t), \\ x_{t+1} &= \arg \min_{x \in \mathcal{K}} \|x - y_{t+1}\|_{A_t}^2. \end{aligned}$$

Remark 10. *There exists some resemblance with Adagrad in the sense that it can be seen as an agile OMD with adaptive regularization function $R_t(x) = \frac{1}{2}\|x - x_1\|_{A_t}$. A major difference is that the diagonal of A_t in the regularization function R_t are equal to the square of the weights in Adagrad, namely D_t^2 , of the form*

$$\frac{1}{t} \frac{1}{t^{-1} \sum_{k=1}^t \nabla f_k(x_k)_i^2}.$$

Theorem 11. *ONS with f_t μ -exp-concave and $\gamma = \frac{1}{2} \min\{\frac{1}{4GD}, \mu\}$ and $\epsilon = 1/(\gamma D)^2$ achieves a regret*

$$\text{Regret}_T \leq 2\left(\frac{1}{\mu} + 4GD\right)d \log T, \quad T \geq 3.$$

Proof. We improve the gradient trick using Proposition 6

$$\sum_{t=1}^T f_t(x_t) - f_t(x^*) \leq \sum_{t=1}^T \nabla f_t(x_t)^T (x_t - x^*) - \frac{\gamma}{2} \sum_{t=1}^T \|x_t - x^*\|_{\nabla f_t(x_t) \nabla f_t(x_t)^T}^2.$$

Using the recursion and the Pythagorean theorem (still valid) we get

$$\begin{aligned} \|x_{t+1} - x^*\|_{A_t}^2 &\leq \|y_{t+1} - x^*\|_{A_t}^2 \\ &\leq (y_{t+1} - x^*)^T A_t (y_{t+1} - x^*) \\ &\leq \|x_t - x^*\|_{A_t}^2 - \frac{1}{\gamma^2} \|\nabla f_t(x_t)\|_{A_t^{-1}}^2 - \frac{2}{\gamma} \nabla f_t(x_t)^T (x_t - x^*). \end{aligned}$$

Thus we get

$$\begin{aligned} \sum_{t=1}^T \nabla f_t(x_t)^T (x_t - x^*) &\leq \frac{\gamma}{2} \sum_{t=1}^T (\|x_t - x^*\|_{A_t}^2 - \|x_{t+1} - x^*\|_{A_t}^2) + \frac{2}{\gamma} \sum_{t=1}^T \|\nabla f_t(x_t)\|_{A_t^{-1}}^2 \\ &\leq \frac{\gamma}{2} \left(\sum_{t=2}^T (\|x_t - x^*\|_{A_t}^2 - \|x_t - x^*\|_{A_{t-1}}^2) + \|x_1 - x^*\|_{A_1}^2 \right) \\ &\quad + \frac{1}{2\gamma} \sum_{t=1}^T \|\nabla f_t(x_t)\|_{A_t^{-1}}^2 \\ &\leq \frac{\gamma}{2} \left(\sum_{t=2}^T \|x_t - x^*\|_{\nabla f_t(x_t) \nabla f_t(x_t)^T}^2 + \|x_1 - x^*\|_{A_1}^2 \right) \\ &\quad + \frac{1}{2\gamma} \sum_{t=1}^T \|\nabla f_t(x_t)\|_{A_t^{-1}}^2 \end{aligned}$$

We immediately derive

$$\text{Regret}_T \leq \frac{\gamma}{2} \|x_1 - x^*\|_{A_1 - \nabla f_1(x_1) \nabla f_1(x_1)^T}^2 + \frac{1}{2\gamma} \sum_{t=1}^T \|\nabla f_t(x_t)\|_{A_t^{-1}}^2.$$

The first term in the upper bound is equal to $\epsilon \|x_1 - x^*\|^2 \leq 1/\gamma^2$. We upper-bound the second term such as

$$\begin{aligned} \sum_{t=1}^T \|\nabla f_t(x_t)\|_{A_t^{-1}}^2 &= \sum_{t=1}^T \text{Tr}(A_t^{-1} \nabla f_t(x_t) \nabla f_t(x_t)^T) \\ &\leq \sum_{t=1}^T \text{Tr}(A_t^{-1} (A_t - A_{t-1})) \\ &\leq \sum_{t=1}^T \log(|A_t|/|A_{t-1}|) \\ &\leq \log(|A_T|/|A_0|). \end{aligned}$$

Since $A_T = \sum_{t=1}^T \nabla f_t(x_t) \nabla f_t(x_t)^T + \epsilon I_d$ then $|A_T| \leq (TG^2 + \epsilon)^d$ and

$$|A_T|/|A_0| \leq (1 + TG^2/\epsilon)^d \leq (1 + TG^2\gamma^2 D^2)^d \leq (1 + T/8)^d \leq T^d$$

for any $T \geq 2$. We get the bound

$$\text{Regret}_T \leq \frac{1}{2\gamma}(1 + d \log T) \leq (4GD + 1/\mu)(1 + d \log T) \leq 2(4GD + 1/\mu)d \log T$$

since $1/\gamma = 2 \max(4GD, 1/\mu) \leq 2(4GD + 1/\mu)$ and $T \geq 3$. \square

Each recursion of the ONS would require to invert a large $\times d$ matrix A_t . Actually one should avoid such inversion by considering the Sherman-Morrisson formula which provides the recursion on A_t^{-1}

$$A_t^{-1} = (A_{t-1} + \nabla f_t(x_t) \nabla f_t(x_t)^T)^{-1} = A_{t-1}^{-1} - \frac{A_{t-1}^{-1} \nabla f_t(x_t) \nabla f_t(x_t)^T A_{t-1}^{-1}}{1 + \nabla f_t(x_t)^T A_{t-1}^{-1} \nabla f_t(x_t)}.$$

Thus the recursion on A_t should be accompanied with one on A_t^{-1} . Moreover the projection $\arg \min_{x \in \mathcal{K}} \|x - y_{t+1}\|_{A_t}^2$ is not explicit for A_t non diagonal (up to my knowledge). One could approximate it with $\arg \min_{x \in \mathcal{K}} \|x - y_{t+1}\|_{\text{Diag}(A_t)}^2$ where $\text{Diag}(A_t)$ is the diagonal matrix extracted from A_t . Then the total cost of one recursion is $O(d^2)$.

One can implement the ONS on MNIST

Algorithm 21: ONS for linear SVM

Parameters: Epoch T , radius $z > 0$, regularization parameter $\lambda > 0$ and $\gamma > 0$.

Initialization: Initial point $x_1 = y_1 = 0$, $A_0 = 1/\gamma^2 I_d$ and $A_0^{-1} = \gamma^2 I_d$.

Sample uniformly iid: $(I_t)_{1 \leq t \leq T}$ from $\{1 \leq i \leq n\}$

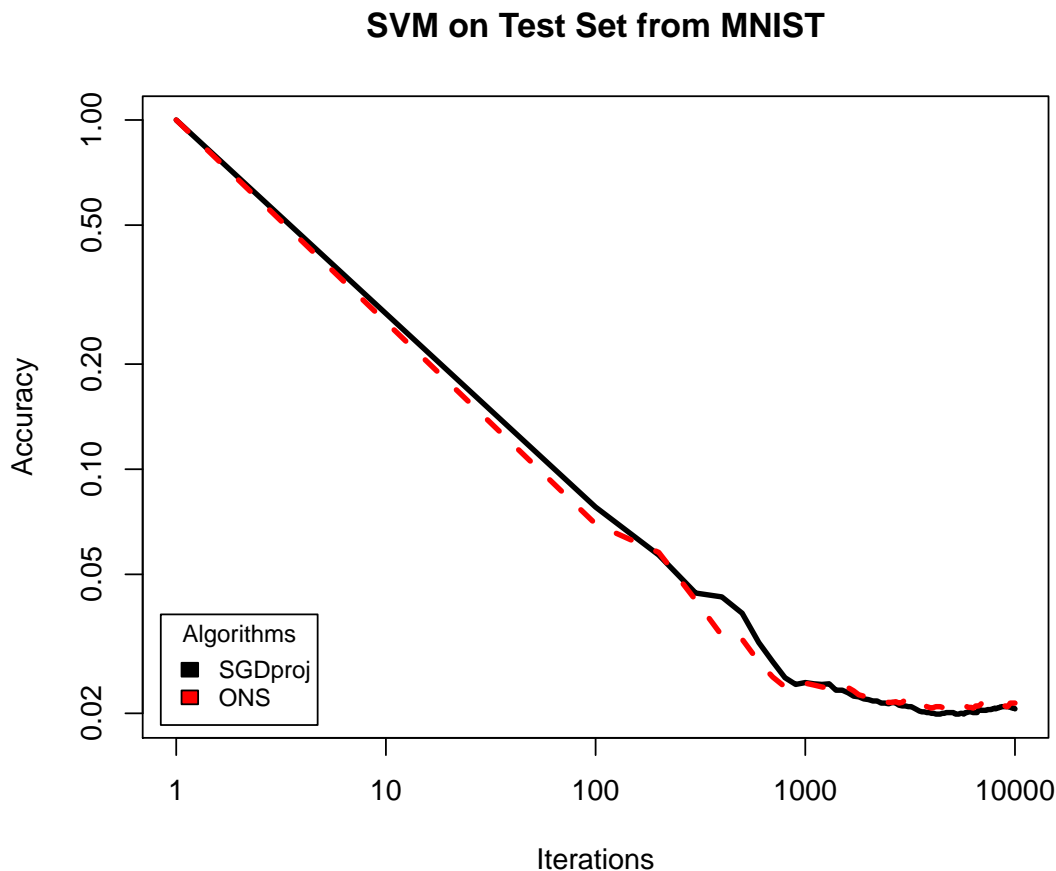
For each iteration $t = 1, \dots, T$:

Iteration: Update

$$\begin{aligned} \nabla_t &= \nabla \ell_{a_t, b_t}(x_t) + \lambda x_t \\ A_t &= A_{t-1} + \nabla_t \nabla_t^T \\ A_t^{-1} &= A_{t-1}^{-1} - \frac{A_{t-1}^{-1} \nabla_t \nabla_t^T A_{t-1}^{-1}}{1 + \nabla_t^T A_{t-1}^{-1} \nabla_t} \\ y_{t+1} &= x_t - \frac{1}{\gamma} A_t^{-1} \nabla_t, \\ x_{t+1} &= \arg \min_{x \in B_1(z)} \|x - y_{t+1}\|_{A_t}^2, \quad 1 \leq i \leq d. \end{aligned}$$

Return: $\bar{x}_{T+1} = \frac{1}{T+1} \sum_{t=1}^{T+1} x_t$

It slightly improves SGD but at the price of a recursion step at $O(d^2)$. Actually it is very complicated to calibrate and the choice of γ is tricky. The performances are very similar to the regularized SGD, both with regularization parameter taken as $\lambda = 1/3$. However the loss in speed is a relative factor of 20.



4.2.3 Natural gradient and EKF

For statistical settings such as binary classification (a_t, b_t) iid, ONS might be useful for statistical purposes. Indeed one updates a matrix A_t^{-1} such that

$$\frac{1}{T} \sum_{t=1}^T A_t^{-1} \approx \mathbb{E}[\nabla \ell_{(a,b)}(x^*) \nabla \ell_{(a,b)}(x^*)^T]^{-1}$$

for T sufficiently large so that $\nabla f_T(x_T) \approx \nabla f_T(x^*)$. One recognizes the variance of the vector score associated to a model

$$(a_t, b_t) \sim c \exp(-\ell_{a,b}(x^*)),$$

where $c > 0$ is the normalizing constant $c = \int_{\mathbb{R}^d} \exp(-\ell_{a,1}(x^*)) da + \int_{\mathbb{R}^d} \exp(-\ell_{a,-1}(x^*)) da$. For regular models we can have the identity

$$\mathbb{E}[\nabla \ell_{(a,b)}(x^*) \nabla \ell_{(a,b)}(x^*)^T] = \mathbb{E}[\nabla^2 \ell_{(a,b)}(x^*)].$$

In such a case \bar{x}_{T+1} might be a good approximation of the maximum likelihood estimator and $\mathbb{E}[\nabla \ell_{(a,b)}(x^*) \nabla \ell_{(a,b)}(x^*)^T]^{-1}$ its associated asymptotic variance

$$\sqrt{T}(\bar{x}_{T+1} - x^*) \sim \mathcal{N}(0, \mathbb{E}[\nabla \ell_{(a,b)}(x^*) \nabla \ell_{(a,b)}(x^*)^T]^{-1}).$$

That ONS provides an estimator of the asymptotic variance $\frac{1}{T} \sum_{t=1}^T A_t^{-1}$ is very useful for instance for significancy testing.

The Online Natural Gradient (or Stochastic Newton) approach replaces $\nabla_t \nabla_t^T$ in the recursion on A_t^{-1} with the better (optimistic) approximation of the variance of the score

$$\mathbb{E}_{(a,b) \sim c \exp(-\ell_{a,b}(x_t))} [\nabla \ell_{(a,b)}(x_t) \nabla \ell_{(a,b)}(x_t)^T].$$

Such quantity is explicit in exponential family under its natural parametrization and then the Stochastic Newton algorithm coincides with the static Extended Kalman Filter. Thus we are forced to consider the logistic model associated with the loss function

$$\frac{(b+1)x^T a}{2} - \log(1 + e^{x^T a}).$$

Algorithm 22: EKF for linear SVM, Fahrmeir (1992)

Parameters: Epoch T .

Initialization: Initial point $x_1 = 0$ and $P_0 = I_d$.

Sample uniformly iid: $(I_t)_{1 \leq t \leq T}$ from $\{1 \leq i \leq n\}$

For each iteration $t = 1, \dots, T$:

Iteration: Update

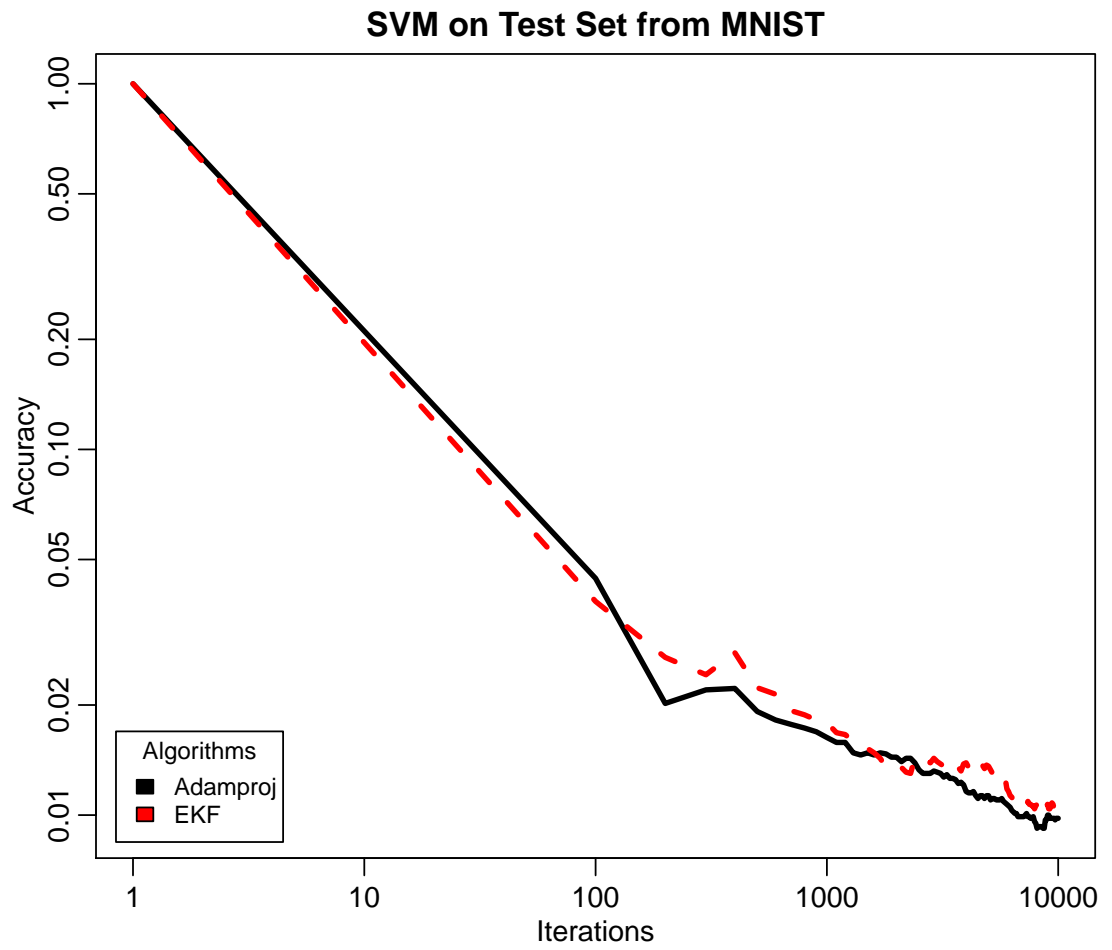
$$\hat{p}_t = \frac{e^{b_{I_t} x_t^T a_{I_t}}}{1 + e^{b_{I_t} x_t^T a_{I_t}}}$$

$$P_t = P_{t-1} - \frac{P_{t-1} a_{I_t} a_{I_t}^T P_{t-1}}{1/(\hat{p}_t(1 - \hat{p}_t)) + a_{I_t}^T P_{t-1} a_{I_t}}$$

$$x_{t+1} = x_t + b_{I_t}(1 - \hat{p}_t)P_t a_{I_t}.$$

Return: x_{T+1}

Note that the recursion depends heavily on the parameter \hat{p}_t which is the sigmoid function applied to $b_{I_t} x_t^T a_{I_t}$ and corresponds to the probability of observing b_{I_t} from a_{I_t} in the logistic model driven by x_t . Thus one can also EKF as an explicit version of a Recursive Bayesian algorithm. In such approach there is no need of the projection nor the averaging step. The obtained accuracy is very close to the one of Adam and the relative loss of speed is only of a relative factor of 10. The gain of relative speed of 1/2 compared with ONS is due to the use of a unique automatically fine tune matrix P rather than the use of two more degenerate matrices A and A^{-1} . Optimal regret bounds have been derived for EKF in such stochastic setting in De Vilmarest and Wintenberger (2021).



Chapter 5

Exploration

5.1 Bandit Convex Optimization

When the information is not complete, we have a new setting called the bandit setting. For instance, in extremely high dimension, even the inquiry of one gradient $\nabla f_t(x_t) \in \mathbb{R}^d$ might be too costly. Instead we use only one coordinate of the gradient $\nabla f_t(x_t)_{I_t} \in \mathbb{R}$ at each step. The information about the gradient is incomplete and the problem is called a Bandit Convex Optimization (BCO) problem.

Remark 11. *This definition of Bandit Convex Optimization departs from the one of Hazan (2019).*

A general reduction from BCO setting to the OCO setting for any \mathcal{K} is to replace $\nabla f_t(x_t)$ by the basic unbiased estimator $\widehat{\nabla f_t(x_t)} = d \nabla f_t(x_t)_{I_t} e_{I_t}$ where $\{e_i\}$ are the element of the canonical basis and I_t are iid uniform over $1, \dots, d$. Indeed then

$$\mathbb{E}[\widehat{\nabla f_t(x_t)}] = \sum_{j=1}^d d \nabla f_t(x_t)_j e_j \mathbb{P}(I_t = j) = \sum_{j=1}^d \nabla f_t(x_t)_j e_j = \nabla f_t(x_t).$$

Thus BCO algorithms are given as stochastic OCO algorithms that explores randomly the space at each recursion. Note that the rest of the OCO algorithm remains unchanged; for instance one can think of an OMD in a randomized version predicting $\widehat{\nabla f_t(x_t)}$ thanks to one coordinate

$$\theta_{t+1} = \theta_t - \eta \widehat{\nabla f_t(x_t)}, \quad x_{t+1} = \nabla R^*(\theta_{t+1}).$$

It is important to notice that for any norm $\|\cdot\|_t^*$ so that $\|e_j\|_t^* = 1$ one has

$$\mathbb{E} \left[\|\widehat{\nabla f_t(x_t)}\|_t^{*2} \right] = \sum_{j=1}^d \|d \nabla f_t(x_t)_j e_j\|_t^{*2} \mathbb{P}(I_t = j) = d \sum_{j=1}^d \nabla f_t(x_t)_j^2 = dG^2,$$

and it is independent of the dual norm. The OCO regret bounds obtained above for OMD methods turn into a BCO regret bound of the form

$$\begin{aligned} \mathbb{E}[\text{Regret}_T] &\leq \mathbb{E} \left[\frac{\eta \sum_{t=1}^T \|\widehat{\nabla f_t(x_t)}\|_t^{*2}}{2} + \frac{D_R^2}{\eta} \right] \\ &\leq \frac{\eta d T G^2}{2} + \frac{D_R^2}{\eta}. \end{aligned}$$

Optimizing in $\eta = D_R/G_{R^*}\sqrt{2/(dT)}$, the randomized algorithm achieves a regret bound deteriorated from the OCO setting by a factor at least \sqrt{d} . For instance, one can randomized the EG+/- and we get a regret bound on the unit ℓ^1 -ball as

$$\mathbb{E}[\text{Regret}_T] \leq zG\sqrt{2dT \log d}$$

by identifying $D_R = \log(d)$. The online to batch conversion still holds on the stochastic version of this algorithm called SREG+/- for which we get

$$\mathbb{E}[h_T^R] \leq \frac{zdG_\infty\sqrt{2\log d}}{\sqrt{T}}$$

using $G \leq \sqrt{d}G_\infty$ and the loss is a factor d . It is worth noticing that the exploration and exploitation are totally independent as the algorithm keeps exploring new directions randomly at each time step.

Algorithm 23: SREG+/- for linear SVM

Parameters: Epoch T , radius $z > 0$.

Initialization: Initial point $x_1 = 0$, weights $w_1 = 1/(2d) \mathbb{1}$.

Sample uniformly iid: $(I_t)_{1 \leq t \leq T}$ from $\{1 \leq i \leq n\}$

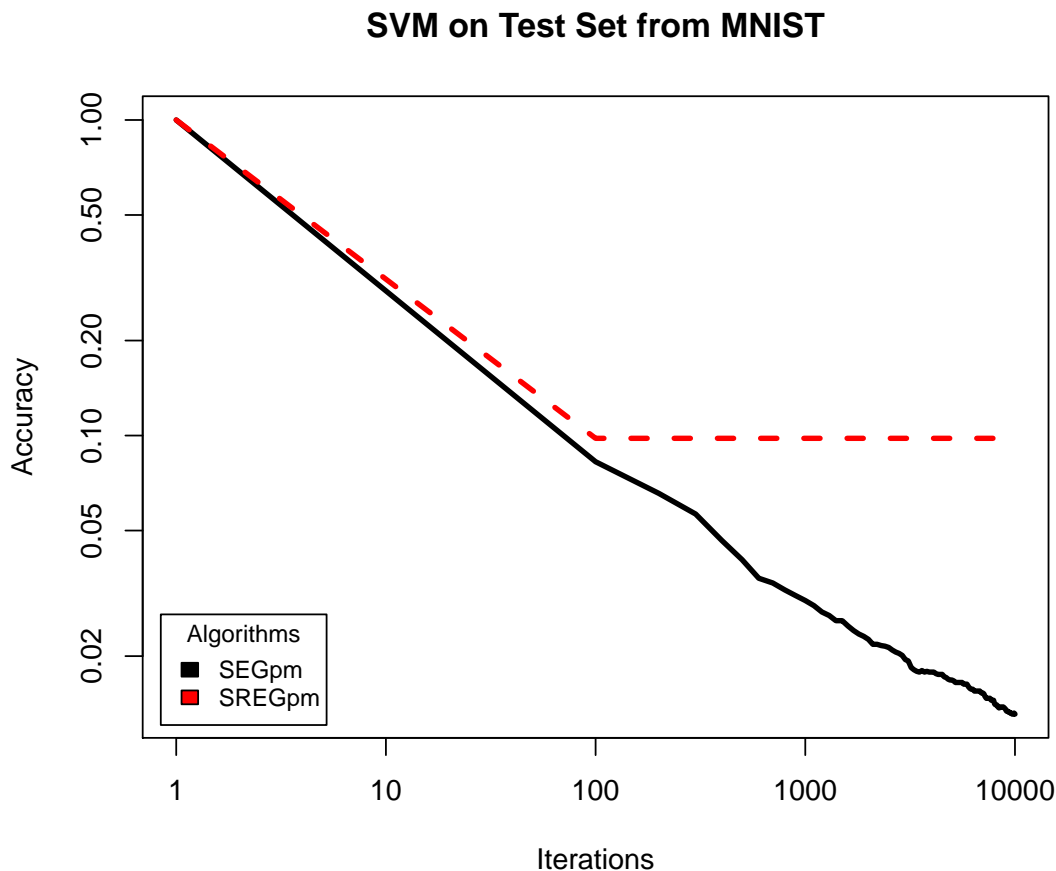
For each iteration $t = 1, \dots, T$:

Sample a direction: $J_t \in \{1, \dots, d\}$ uniformly

Iteration: Update

$$\begin{aligned} \eta_t &= 1/\sqrt{dt} \\ w_{t,J_t} &\leftarrow \exp(-\eta_t d \nabla \ell_{a_{I_t}, b_{I_t}}(x_t)_{J_t}) w_{t,J_t}, \\ w_{t,J_t+d} &\leftarrow \exp(\eta_t d \nabla \ell_{a_{I_t}, b_{I_t}}(x_t)_{J_t}) w_{t,J_t+d}, \\ w_{t+1} &= \frac{w_t}{\sum_{i=1}^{2d} w_{t,i}} \\ x_{t+1,i} &= z(w_{t+1,i} - w_{t+1,i+d}), \quad 1 \leq i \leq d. \end{aligned}$$

Return: $\bar{x}_{T+1} = \frac{1}{T+1} \sum_{t=1}^{T+1} x_t$



The algorithm is stuck at the accuracy 0.1, the accuracy of the random guess (as there is only 10% of label 1 in the dataset, digit 0). The algorithm explores at each round and does not achieve a good exploration-exploitation trade-off.

5.2 Exp3 algorithm

Similarly as in the Expert Advice setting, at each round the algorithm assigns confident weights $x_t \in \Lambda$, pick randomly an expert $I_t \sim x_t$ and incur the loss ℓ_{t,I_t} . However the information of the algorithm is now limited to the loss ℓ_{t,I_t} but the regret is unchanged

$$\mathbb{E}[\text{Regret}_T(\ell)] = \sum_{t=1}^T \mathbb{E}_{x_t}[\ell_t] - \min_{1 \leq i \leq d} \sum_{t=1}^T \ell_{t,i}.$$

At each round the algorithm can either explore and pick a new expert that has been never played or exploit and pick an already chosen expert in order to learn its performances. This is an exploration-exploitation trade-off called the Multi-Armed Bandit (MAB) problem. Note that in the bandit setting the experts are also called actions since one could also argue that at each step only one expert is active.

A solution for obtaining the trade-off exploration-exploitation has been provided by copying the EWA replacing ℓ_t in the exponential by some unbiased estimator. Considering

$\ell_{t,i}$ for the coordinate which coincides with the chosen one $i = I_t$ and 0 elsewhere, we get

$$\mathbb{E}_{x_t}[\ell_{t,i} \mathbb{1}_{i=I_t}] = \sum_{j=1}^d \ell_{t,i} \mathbb{1}_{i=j} \mathbb{P}(I_t = j) = \ell_{t,i} x_{t,i}, \quad 1 \leq i \leq d,$$

which is a biased estimation of $\ell_{t,i}$. Indeed the distribution of I_t makes the strategy biased in favor of the experts with large confidence weights. In order to debiased the approximation we consider instead

$$\widehat{\ell}_{t,i} = \frac{\ell_{t,i}}{x_{t,i}} \mathbb{1}_{i=I_t}.$$

Then we get an unbiased estimator of $\ell_{t,i}$, indeed

$$\mathbb{E}_{x_t}[\widehat{\ell}_{t,i}] = \sum_{j=1}^d \frac{\ell_{t,i}}{x_{t,i}} \mathbb{1}_{i=j} \mathbb{P}(I_t = j) = \sum_{j=1}^d \frac{\ell_{t,i}}{x_{t,i}} \mathbb{1}_{i=j} x_{t,j} = \ell_{t,i}, \quad 1 \leq i \leq d.$$

We obtain the Exp3 (Exponential weights for Exploration and Exploitation) algorithm

Algorithm 24: Exp3 algorithm (simple version), Auer et al. (2002)

Parameters: step-size $\eta > 0$.

Initialization: Initial prediction $x_1 = (1/d) \mathbb{1}$.

For each recursion $t \geq 1$:

Sample an expert: $I_t \sim x_t$ uniformly

Predict as the I_t -th expert

Incur the average loss: ℓ_{t,I_t}

Observe: $\ell_{t,I_t} \in \mathbb{R}$

Recursion: Update

$$\widehat{\ell}_{t,i} = \frac{\ell_{t,i}}{x_{t,i}} \mathbb{1}_{i=I_t}, \quad 1 \leq i \leq d$$

$$x_{t+1} = \frac{\exp(-\eta \widehat{\ell}_t) x_t}{\sum_{i=1}^d \exp(-\eta \widehat{\ell}_{t,i}) x_{t,i}}.$$

We obtain the regret bound in the case of non-negative losses:

Theorem 12. *In the MAB setting with $\|\ell_t\| \leq G$ and $\ell_t \geq 0$, for $\eta = G^{-1} \sqrt{2 \log d/T}$ we obtain a regret bound for Exp3 as*

$$\mathbb{E}[\text{Regret}_T(\ell)] \leq G \sqrt{2T \log d}.$$

Proof. We refine the previous analysis of EWA for random loss $\widehat{\ell}_t > 0$ and $\eta > 0$. Recall that now the Exp3 strategy x_t itself is random (it depends on the past sampled coordinates I_s , $s < t$) in order to get the regret bound

$$\mathbb{E}[\text{Regret}_T(\ell)] = \mathbb{E}[\text{Regret}_T(f)] = \mathbb{E} \left[\sum_{t=1}^T x_t^T \ell_t - \min_{x \in \Lambda} \sum_{t=1}^T x^T \ell_t \right]$$

for the linear loss $f_t(x_t) = x_t^T \ell_t$. Because $\widehat{\nabla} f_t(x_t) = \widehat{\ell}_t = \ell_{t,I_t}/x_{t,I_t} e_{I_t}$ is an unbiased estimator of ℓ_t we also have

$$\mathbb{E} \left[\sum_{t=1}^T x_t^T \ell_t - \min_{x \in \Lambda} \sum_{t=1}^T x^T \ell_t \right] = \mathbb{E} \left[\sum_{t=1}^T x_t^T \widehat{\ell}_t - \min_{x \in \Lambda} \sum_{t=1}^T x^T \widehat{\ell}_t \right] = \mathbb{E}[\text{Regret}_T(\widehat{\nabla} f)].$$

We notice that Exp3 is an OMD algorithm using a stochastic unbiased approximation of $\widehat{\nabla f_t(x_t)} = \widehat{\ell}_t = \ell_{t,I_t}/x_{t,I_t}e_{I_t}$. Thus we obtain an expected regret bounded of the form

$$\mathbb{E}[\text{Regret}_T(\widehat{\nabla f})] \leq \mathbb{E}\left[\frac{\sum_{t=1}^T \|\eta \widehat{\nabla f_t(x_t)}\|_t^{*2}}{2\eta} + \frac{\log d}{\eta}\right].$$

We cannot use anymore the rough bound $\|\cdot\|_t^* \leq \|\cdot\|_\infty$ because $\|\widehat{\nabla f_t(x_t)}\|_\infty = \ell_{t,I_t}/x_{t,I_t}$ is not bounded since x_{t,I_t} can be as close to 0 as possible. Instead, we refine our bound on the dual norm going back to its definition

$$\frac{1}{2}\|\eta \widehat{\nabla f_t(x_t)}\|_t^{*2} = B_{R^*}(\theta_t - \eta \widehat{\nabla f_t(x_t)}, \theta_t),$$

for $R(x) = x^T \log(x)$ and $\theta_t = \nabla R(y_t)$. Then we make the important remark that the update of Exp3 (and also EWA) can be written in the agile version as in Exercize 11, namely

$$\nabla R(y_{t+1}) = \nabla R(x_t) - \eta \nabla f_t(x_t) \text{ and } x_{t+1} = \arg \min_{x \in \mathcal{K}} B_R(x || y_{t+1}).$$

Thus one can replace θ_t with $\nabla R(x_t)$. As in our setting we have $\widehat{\nabla f_t(x_t)} = \widehat{\ell}_{t,i}$, we get

$$\begin{aligned} B_{R^*}(\theta_t - \eta \widehat{\ell}_{t,i}, \theta_t) &= B_{R^*}(\nabla R(x_t) - \eta \widehat{\ell}_{t,i}, \nabla R(x_t)) \\ &= R^*(\nabla R(x_t) - \eta \widehat{\ell}_{t,i}) - R^*(\nabla R(x_t)) \\ &\quad + \eta \nabla R^*(\nabla R(x_t))^T \widehat{\ell}_{t,i}. \end{aligned}$$

Since $R^*(x^*) = y^{*T} x^* - R(y^*)$ with $\nabla R(y^*) = x^*$, we get

$$\begin{aligned} R^*(\nabla R(x_t)) &= x_t^T \nabla R(x_t) - R(x_t) = x_t^T \mathbb{1}, \\ R^*(\nabla R(x_t) - \eta \widehat{\ell}_{t,i}) &= (x_t e^{-\eta \widehat{\ell}_{t,i}})^T (\nabla R(x_t) - \eta \widehat{\ell}_{t,i}) - R(x_t e^{-\eta \widehat{\ell}_{t,i}}) \\ &= (x_t e^{-\eta \widehat{\ell}_{t,i}})^T \mathbb{1}. \end{aligned}$$

Finally, using the relation $\nabla R^*(\nabla R(x_t)) = x_t$ we get

$$B_{R^*}(\theta_t - \eta \widehat{\ell}_{t,i}, \theta_t) = x_t^T (e^{-\eta \widehat{\ell}_{t,i}} - 1 + \eta \widehat{\ell}_{t,i}).$$

Since $\exp(-x) - x - 1 \leq x^2/2$ for any $x > 0$ we get the desired improved regret bound

$$\text{Regret}_T(\widehat{\nabla f}) \leq \frac{\eta}{2} \sum_{t=1}^T x_t^T \widehat{\ell}_t^2 + \frac{\log d}{\eta}.$$

Note that the regret bound is still depending on the randomness of i_t via $\widehat{\ell}_t$. We take its expectation

$$\mathbb{E}[\text{Regret}_T(\ell)] \leq \mathbb{E}\left[\frac{\eta}{2} \sum_{t=1}^T x_t^T \mathbb{E}_{x_t}[\widehat{\ell}_t^2] + \frac{\log d}{\eta}\right].$$

Thus we have to upper bound

$$\mathbb{E}_{x_i}[\widehat{\ell}_{t,i}^2] = \sum_{j=1}^d \left(\frac{\ell_{t,i}}{x_{t,i}}\right)^2 \mathbb{1}_{i=j} \mathbb{P}(I_t = j) = \frac{\ell_{t,i}^2}{x_{t,i}}$$

and we obtain

$$\mathbb{E}[\text{Regret}_T(\ell)] \leq \frac{\eta}{2} \sum_{t=1}^T \sum_{i=1}^d \ell_{t,i}^2 + \frac{\log d}{\eta} \leq \frac{\eta}{2} TG^2 + \frac{\log d}{\eta},$$

and the desired result follows. \square

The regret bound is less accurate than in the complete information setting with a relative loss of G that can be as large as \sqrt{d} . It is still optimal up to log terms in this incomplete information setting since one has to explore the space. Exp3 can be improved as the unbiased estimator $\widehat{\nabla f_t(x_t)}$ is not bounded since the confident weights $x_{t,i}$ can be as small as possible. Note that on the contrary to the EWA algorithm, the recursion is not invariant by a shift of the loss and the condition of non-negativity of the losses is necessary.

5.3 Exp2 algorithm for OCO on $\mathcal{K} = B_1(z)$

Another reduction of BCO to the OCO setting when $\mathcal{K} = B_1(z)$ is the following one. We first apply the gradient trick in order to minimize the linearized regret

$$\begin{aligned} \mathbb{E}[\text{Regret}_T(f)] &= \mathbb{E} \left[\sum_{t=1}^T f_t(x_t) - \min_{x \in \Lambda} \sum_{t=1}^T f_t(x) \right] \\ &\leq \mathbb{E} \left[\sum_{t=1}^T \nabla f_t(x_t)^T (x_t - x^*) \right] = \mathbb{E}[\text{Regret}_T(\nabla f)] \end{aligned}$$

by considering the linearized losses $f_t(x) = \nabla f_t(x_t)^T x$, $x \in B_1(z)$. One has to consider a set of actions as $A_t \in \{e_1, \dots, e_{2d}\}$ associated to confidence weights $w \in \Lambda_{2d}$ and consider the stochastic approximation

$$\pm \widehat{\nabla f_t(x_t)} = \frac{\pm \nabla f_t(x_t)^T A_t}{w_{A_t}} A_t,$$

where w_{A_t} is the confidence weights assigned to the action A_t . The approximation of the signed gradient is still unbiased

$$\mathbb{E}[\pm \widehat{\nabla f_t(x_t)}_i] = \mathbb{E}_w[\pm \widehat{\nabla f_t(x_t)}_i] = \sum_{j=1}^{2d} \frac{\pm \nabla f_t(x_t)_i}{w_i} \mathbb{1}_{i=j} \mathbb{P}(A_t = j) = \pm \nabla f_t(x_t)_i, \quad 1 \leq i \leq 2d.$$

The weights are updated as in EWA on $2d$ experts and the prediction x_t is provided by $x_{t,i} = z(w_{t,i} - w_{t,i+d})$ as in EG+/- . The same analysis as for Exp3 is conducted except that the assumption of non-negativeness on the loss does not hold. This issue is dropped by replacing the inequality $\exp(-x) - x - 1 \leq x^2/2$ only valid for $x > 0$ by the inequality $\exp(-x) - x - 1 \leq x^2$ valid for any $|x| \leq 1$ for a learning rate satisfying $\|\eta \widehat{\nabla f_t(x_t)}\|_\infty \leq 1$. In order to do so it is necessary to bound $\widehat{\nabla f_t(x_t)}$ by another trick, the introduction of an exploration factor γ so that the weights are lower bounded.

Algorithm 25: Exp2 algorithm, Bubeck and Cesa-Bianchi (2012) (also called Exp3 for linear bandit, Lattimore and Szepesvári (2020))

Parameters: step-size $\eta > 0$.

Initialization: Initial prediction $w'_1 = w_1 = (1/2d)\mathbb{1}$ and $x_1 = 0$. Exploration rate $\gamma = 2\eta zdG_\infty$.

For each recursion $t \geq 1$:

Sample an action: $A_t \in \{e_1, \dots, e_{2d}\} \sim w_t$ that determines a coordinate of $\pm \widehat{\nabla f_t}(x_t)$

Predict x_t

Incur the average loss: $f_t(x_t)$

Observe: $\nabla f_t(x_t)_{I_t} \in \mathbb{R}$

Recursion: Update

$$w'_{t+1} = \frac{\exp(-\eta \pm \widehat{\nabla f_t}(x_t)) w'_t}{\sum_{i=1}^{2d} \exp(-\eta \pm \widehat{\nabla f_t}(x_t)_i) w'_{t,i}},$$

$$w_{t+1} = (1 - \gamma)w'_{t+1} + \gamma w_1,$$

$$x_{t+1} = z(w'_{t+1,i} - w'_{t+1,i+d})_{1 \leq i \leq d}.$$

Theorem 13. *In the BCO setting with $\|\nabla f_t(x)\|_\infty \leq G_\infty$ for any $x \in \mathcal{K}$, choosing $\eta = (zG)^{-1} \sqrt{\log(2d)/(2T)}$ and $\gamma = 2\eta zdG_\infty$ we obtain a regret bound for Exp2 on $B_1(1)$ as*

$$\mathbb{E}[\text{Regret}_T] \leq zG \sqrt{2T \log(2d)}, \quad T > 2d \log(2d) (G_\infty/G)^2.$$

Proof. By definition of the weights and using Hölder inequality, we have

$$\begin{aligned} \mathbb{E}[\text{Regret}_T(\nabla f)] &= \mathbb{E}[\text{Regret}_T(\widehat{\nabla} f)] \\ &\leq \mathbb{E}\left[\sum_{t=1}^T \widehat{\nabla f_t}(x_t)^T (x_t - x^*)\right] \\ &\leq \mathbb{E}\left[\sum_{t=1}^T \widehat{\nabla f_t}(x_t)^T (z(w'_{t+1,i} - w'_{t+1,i+d})_{1 \leq i \leq d} - x^*)\right] \\ &\leq \mathbb{E}\left[\sum_{t=1}^T z \pm \widehat{\nabla f_t}(x_t)^T (w'_t - w^*)\right], \end{aligned}$$

where $w^* \in \Lambda_{2d}$ by an application of Lemma 1. Because $w_t \geq \gamma/(2d)$ and the specific choice of γ we have $\|\eta z \pm \widehat{\nabla f_t}(x_t)\|_\infty \leq 2\eta zdG_\infty/\gamma \leq 1$. We extend to non positive losses (here the estimated gradients) the previous bound on Exp3 at the cost of a factor 2 using $\exp(-x) - x - 1 \leq x^2$ valid for every $|x| \leq 1$. We obtain

$$\mathbb{E}\left[z \sum_{t=1}^T \pm \widehat{\nabla f_t}(x_t)^T (w_t - w^*)\right] \leq \mathbb{E}\left[\eta \sum_{t=1}^T w_t^T \mathbb{E}_{w_t} \left[(z \pm \widehat{\nabla f_t}(w_t))^2\right]\right] + \frac{\log(2d)}{\eta}.$$

We estimate

$$w_t^T \mathbb{E}_{w_t} \left[\pm \widehat{\nabla f_t}(w_t)^2\right] \leq w_t^T \sum_{i=1}^{2d} \left(\frac{\nabla f_t(x_t)_i}{w_{t,i}}\right)^2 w_{t,i} e_i = 2\|\nabla f_t(x_t)\|^2 \leq 2G^2,$$

and we obtain

$$\mathbb{E}[\text{Regret}_T] \leq 2\eta T(zG)^2 + \frac{\log(2d)}{\eta}.$$

We obtain the desired result by the specific choice of η . Note that the restriction on T comes from the fact that γ must be smaller than 1. \square

We provide a version of this algorithm for SVM on MNIST called Stochastic Bandit EG+/- . The online to batch conversion still holds on SBEG+/- for which we get

$$\mathbb{E}[h_T^R] \leq \frac{zG\sqrt{2\log 2d}}{\sqrt{T}}.$$

In theory we gain a factor \sqrt{d} compared with SREG+/- thanks to the exploration that is now adaptive to the OCO problem. However in practice the algorithm is also stuck at the accuracy 0.1, even if the algorithm explores less and less depending on how fast it learns.

Algorithm 26: SBEG+/- for linear SVM

Parameters: Epoch T , radius $z > 0$.

Initialization: Initial point $x_1 = 0$ and weights $w_1 = w'_1 = 1/(2d) \mathbb{1}$

Sample uniformly iid: $(I_t)_{1 \leq t \leq T}$ from $\{1 \leq i \leq n\}$

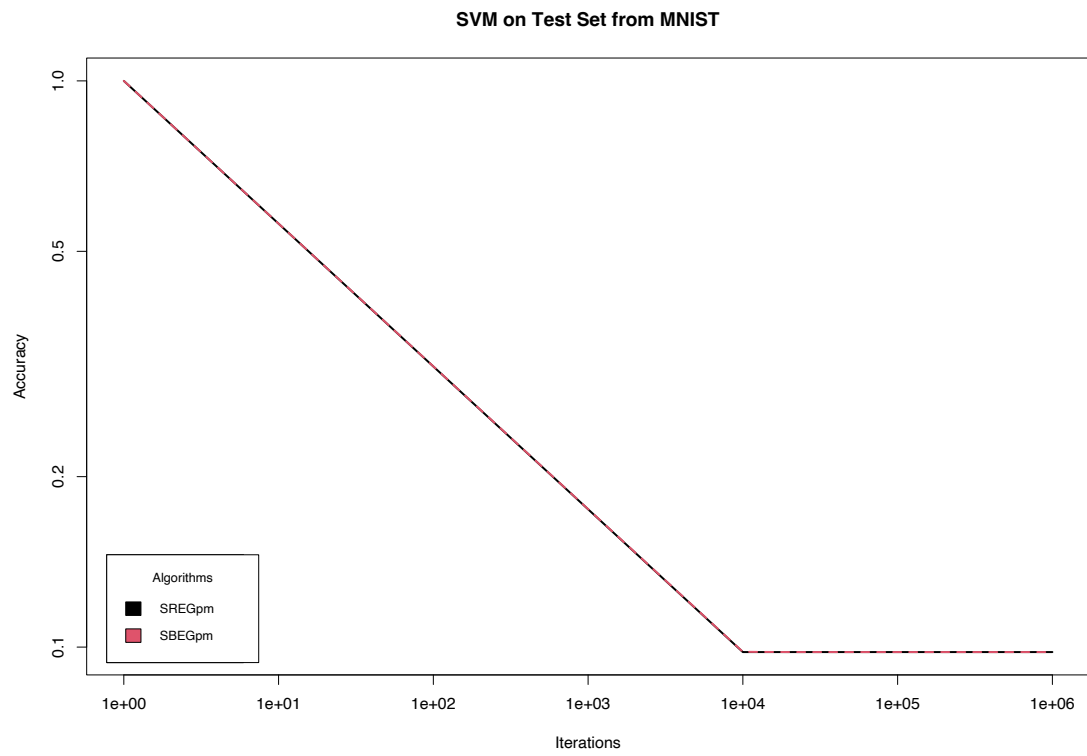
For each iteration $t = 1, \dots, T$:

Sample an action: $A_t \in \{e_1, \dots, e_{2d}\} \sim w_t$ that determines a coordinate of $\pm \nabla \ell_{a_{I_t}, b_{I_t}}(x_t)$

Iteration: Update

$$\begin{aligned} \eta_t &= 1/\sqrt{dt} \\ \gamma_t &= \min(1, d\eta_t) \\ w'_{t,A_t} &\leftarrow \exp(-\eta_t \pm \nabla \ell_{a_{I_t}, b_{I_t}}(x_t)_{A_t}/w_{t,A_t})w'_{t,A_t}, \\ w'_{t+1} &= \frac{w'_t}{\sum_{i=1}^{2d} w'_{t,i}} \\ w_{t+1} &= (1 - \eta_t)w'_t + \frac{\eta_t}{2d} \\ x_{t+1,i} &= z(w'_{t+1,i} - w'_{t+1,i+d}), \quad 1 \leq i \leq d. \end{aligned}$$

Return: $\bar{x}_{T+1} = \frac{1}{T+1} \sum_{t=1}^{T+1} x_t$.



Bibliography

- Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. The nonstochastic multiarmed bandit problem. *SIAM journal on computing*, 32(1):48–77, 2002.
- Stephen Boyd, Stephen P Boyd, and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- Sébastien Bubeck. Convex optimization: Algorithms and complexity. *arXiv preprint arXiv:1405.4980*, 2014.
- Sébastien Bubeck and Nicolo Cesa-Bianchi. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *arXiv preprint arXiv:1204.5721*, 2012.
- Joseph De Vilmorest and Olivier Wintenberger. Stochastic online optimization using kalman recursion. *JMLR*, 2021.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- Ludwig Fahrmeir. Posterior mode estimation by extended kalman filtering for multivariate dynamic generalized linear models. *Journal of the American Statistical Association*, 87(418):501–509, 1992.
- Steffen Grünewälder. Compact convex projections. *The Journal of Machine Learning Research*, 18(1):8089–8131, 2017.
- Elad Hazan. Introduction to online convex optimization. *arXiv preprint arXiv:1909.05207*, 2019.
- Elad Hazan and Satyen Kale. Extracting certainty from uncertainty: Regret bounded by variation in costs. *Machine learning*, 80(2-3):165–188, 2010.
- Elad Hazan and Satyen Kale. Beyond the regret minimization barrier: an optimal algorithm for stochastic strongly-convex optimization. In *Proceedings of the 24th Annual Conference on Learning Theory*, pages 421–436, 2011.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Jyrki Kivinen and Manfred K Warmuth. Exponentiated gradient versus gradient descent for linear predictors. *information and computation*, 132(1):1–63, 1997.

- Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- Nick Littlestone and Manfred K Warmuth. The weighted majority algorithm. *Information and computation*, 108(2):212–261, 1994.
- Yurii E Nesterov. A method for solving the convex programming problem with convergence rate $O(1/k^2)$. In *Dokl. akad. nauk Sssr*, volume 269, pages 543–547, 1983.
- Francesco Orabona. A modern introduction to online learning. *arXiv preprint arXiv:1912.13213*, 2019.
- Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. *arXiv preprint arXiv:1904.09237*, 2019.
- Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.
- Shai Shalev-Shwartz and Yoram Singer. A primal-dual perspective of online learning algorithms. *Machine Learning*, 69(2-3):115–142, 2007.
- Shai Shalev-Shwartz et al. Online learning and online convex optimization. *Foundations and trends in Machine Learning*, 4(2):107–194, 2011.
- Olivier Wintenberger. Optimal learning with bernstein online aggregation. *Machine Learning*, 106(1):119–141, 2017.
- Martin Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Proceedings of the 20th international conference on machine learning (icml-03)*, pages 928–936, 2003.