## Exercise sheet

Lecturer: Georgios P. Karagiannis

georgios.karagiannis@durham.ac.uk

## Part 1. Elements of convex learning problems

**Exercise 1.**  $(\star)$ Let  $f: \mathbb{R}^d \to \mathbb{R}$  such that  $f(w) = g(\langle w, x \rangle + y)$  or some  $x \in \mathbb{R}^d$ ,  $y \in \mathbb{R}$ . Show that: If g is convex function then f is convex function.

**Exercise 2.** (\*)Let functions  $g_1$  be  $\rho_1$ -Lipschitz and  $g_2$  be  $\rho_2$ -Lipschitz. Then, show that, f with  $f(x) = g_1(g_2(x))$  is  $\rho_1\rho_2$ -Lipschitz.

**Exercise 3.**  $(\star)$ Let  $f: \mathbb{R}^d \to \mathbb{R}$  with  $f(w) = g(\langle w, x \rangle + y)$   $x \in \mathbb{R}^d$  and  $y \in \mathbb{R}$ . Let  $g: \mathbb{R} \to \mathbb{R}$  be a  $\beta$ -smooth function. Then show that f is a  $(\beta ||x||^2)$ -smooth.

**Hint:** You may use Cauchy-Schwarz inequality  $\langle y, x \rangle \leq ||y|| \, ||x||$ 

**Exercise 4.** (\*)Show that  $f: S \to \mathbb{R}$  is  $\rho$ -Lipschitz over an open convex set S if and only if for all  $w \in S$  and  $v \in \partial f(w)$  it is  $||v|| \le \rho$ .

**Hint::** You may use Cauchy-Schwarz inequality  $\langle y, x \rangle \leq ||y|| \, ||x||$ 

**Exercise 5.** (\*)Let  $g_1(w), ..., g_r(w)$  be r convex functions, and let  $f(\cdot) = \max_{\forall j} (g_j(\cdot))$ . Show that for some w it is  $\nabla g_k(w) \in \partial f(w)$  where  $k = \arg \max_j (g_j(w))$  is the index of function  $g_j(\cdot)$  presenting the greatest value at w.

**Exercise 6.** (\*)Consider the regression learning problem  $(\mathcal{H}, \mathcal{Z}, \ell)$  with predictor rule  $h(x) = \langle w, x \rangle$  labeled by some unknown parameter  $w \in \mathcal{W}$ , loss function  $\ell(w, (x, y)) = (\langle w, x \rangle - y)^2$ , feature  $x \in \mathcal{X}$ , and target  $y \in \mathbb{R}$ . Let  $\mathcal{W} = \mathcal{X} = \{\omega \in \mathbb{R}^d : |\omega| \leq \rho\}$  for some  $\rho > 0$ .

- (1) Show that the resulting learning problem is Convex-Lipschitz-Bounded learning problem.
- (2) Specify the parameters of Lipschitnzess.

**Exercise 7.**  $(\star)$ If f is  $\lambda$ -strongly convex and u is a minimizer of f then for any w

$$f(w) - f(u) \ge \frac{\lambda}{2} \|w - u\|^2$$

**Hint::** Use the definition, and set  $\alpha \to 0$ .

The following is given as a homework (Formative assessment 1)

**Exercise 8.**  $(\star)$  Let  $f: \mathbb{R}^d \to \mathbb{R}$  be a convex and  $\beta$ -smooth function.

(1) Show that for  $v, w \in \mathbb{R}^d$ 

$$f(v) - f(w) \in \left(\left\langle \nabla f(w), v - w \right\rangle, \left\langle \nabla f(w), v - w \right\rangle + \frac{\beta}{2} \|v - w\|^2 \right)$$

(2) Show that for  $v, w \in \mathbb{R}^d$  such that  $v = w - \frac{1}{\beta} \nabla f(w)$ , it is

$$\frac{1}{2\beta} \left\| \nabla f\left(w\right) \right\|^{2} \le f\left(w\right) - f\left(v\right)$$

(3) Additionally assume that f(x) > 0 for all  $x \in \mathbb{R}^d$ . Show that for  $w \in \mathbb{R}^d$ ,

$$\|\nabla f\left(w\right)\| \le \sqrt{2\beta f\left(w\right)}$$

The following is given as a homework (Formative assessment 1)

**Exercise 9.**  $(\star)$ Let  $f: \mathbb{R}^d \to \mathbb{R}$  be a  $\lambda$ -strongly convex function. Assume that  $w^*$  is a minimizer of f i.e.

$$w^* = \operatorname*{arg\,min}_{w} \left\{ f\left(w\right) \right\}$$

Show that for any  $w \in \mathbb{R}^d$  it holds

$$f(w) - f(w^*) \ge \frac{\lambda}{2} \|w - w^*\|^2$$

**Hint:** Use the definition of  $\lambda$ -strongly convex function, properly rearrange it, and ...

**Exercise 10.** (\*)Show that the function  $J(x;\lambda) = \lambda ||x||^2$  is  $2\lambda$ -strongly convex

## Part 2. Stochastic learning

Exercise 11.  $(\star)$  Assume a Bayesian model

$$\begin{cases} z_i | w & \stackrel{\text{ind}}{\sim} f(z_i | w), \ i = 1, ..., n \\ w & \sim f(w) \end{cases}$$

and consider that our objective is the discovery of MAP estimate  $w^*$  i.e.

$$w^* = \arg\min_{\forall w \in \Theta} \left(-\log\left(L_n\left(w\right)\right) - f\left(w\right)\right) = \arg\min_{\forall w \in \Theta} \left(-\sum_{i=1}^n \log\left(f\left(z_i|w\right)\right) - \log\left(f\left(w\right)\right)\right)$$

by using SGD with update

$$w^{(t+1)} = w^{(t)} + \eta_t \left( \frac{n}{m} \sum_{j \in \mathcal{J}^{(t)}} \nabla_w \log \left( f\left(z_j | w^{(t)}\right) \right) + \nabla_w \log \left( f\left(w^{(t)}\right) \right) \right)$$

for some randomly selected set  $\mathcal{J}^{(t)} \subseteq \{1,...,n\}^m$  of m integers from 1 to n via simple random sampling (SRS) with replacement. Show that

$$\mathbf{E}_{\mathcal{J}^{(t)} \sim \text{simple-random-sampling}} \left( \frac{n}{m} \sum_{j \in \mathcal{J}^{(t)}} \nabla_w \log \left( f\left(z_j | w^{(t)}\right) \right) \right) = \sum_{i=1}^n \nabla_w \log \left( f\left(z_i | w^{(t)}\right) \right)$$

**Exercise 12.** (\*) Let  $\{v_t; t = 1, ..., T\}$  be a sequence of vectors. Consider an algorithm producing  $\{w^{(t)}; t = 1, 2, 3, ...\}$  with

$$w^{(1)} = 0$$
$$w^{(t+1)} = w^{(t)} - \eta v_t$$

Show that

(1) it is

$$\langle w^{(t)} - w^*, v_t \rangle = \frac{1}{2\eta} \left( -\left\| w^{(t+1)} - w^* \right\|^2 + \left\| w^{(t)} - w^* \right\|^2 \right) + \frac{\eta}{2} \left\| v_t \right\|^2$$

(2) it is

$$\sum_{t=1}^{T} \langle w^{(t)} - w^*, v_t \rangle = \frac{1}{2\eta} \sum_{t=1}^{T} \left( -\left\| w^{(t+1)} - w^* \right\|^2 + \left\| w^{(t)} - w^* \right\|^2 \right) + \frac{\eta}{2} \sum_{t=1}^{T} \left\| v_t \right\|^2$$

(3) (continue) it is

$$\sum_{t=1}^{T} \langle w^{(t)} - w^*, v_t \rangle \le \frac{\|w^*\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^{T} \|v_t\|^2$$

**Exercise 13.**  $(\star)$  Let  $\{v_t; t=1,...,T\}$  be a sequence of vectors. Consider an algorithm producing  $\{w^{(t)}; t=1,2,3,...\}$  with

$$w^{(1)} = 0$$

$$w^{\left(t + \frac{1}{2}\right)} = w^{(t)} - \eta v_t$$

$$w^{(t+1)} = \arg\min_{w \in \mathcal{H}} \left( \left\| w - w^{\left(t + \frac{1}{2}\right)} \right\| \right)$$

for t = 1, ..., T.

Hint: You can use the following Lemma

(**Projection Lemma**): Let  $\mathcal{H}$  be a closed convex set and let v be the projection of w onto  $\mathcal{H}$ ,i.e.

$$v = \operatorname*{arg\,min}_{x \in \mathcal{H}} \|x - w\|^2$$

then for every  $u \in \mathcal{H}$  it is

$$||v - u||^2 \le ||w - u||^2$$

Show that

(1) it is

$$\langle w^{(t)} - w^*, v_t \rangle \le \frac{1}{2\eta} \left( -\left\| w^{(t+1)} - w^* \right\|^2 + \left\| w^{(t)} - w^* \right\|^2 \right) + \frac{\eta}{2} \left\| v_t \right\|^2$$

(2) it is

$$\sum_{t=1}^{T} \langle w^{(t)} - w^*, v_t \rangle \le \frac{1}{2\eta} \sum_{t=1}^{T} \left( -\left\| w^{(t+1)} - w^* \right\|^2 + \left\| w^{(t)} - w^* \right\|^2 \right) + \frac{\eta}{2} \sum_{t=1}^{T} \left\| v_t \right\|^2$$

(3) (continue) it is

$$\sum_{t=1}^{T} \langle w^{(t)} - w^*, v_t \rangle \le \frac{\|w^*\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^{T} \|v_t\|^2$$

Comment: Above we show that Lemma ?? from "Handout ??: Gradient descent" holds even when a projection step is included. Hence, even if a projection step is included after the update step of the recursion of GD algorithm or the SGD algorithm the analysis in Section ?? in "Handout ??: Gradient descent" holds. Hence, even if a projection step is included after the update step of the recursion of SGD algorithm or the SGD algorithm the analysis in Section ?? in "Handout ??: Stochastic gradient descent" holds.

The following is given as a homework (Formative assessment 2)

**Exercise 14.** (\*) <sup>1</sup>Consider the binary classification problem with inputs  $x \in \mathcal{X}$  where  $\mathcal{X} := \{x \in \mathbb{R}^d : ||x||_2 \le L\}$  for some given value L > 0, target  $y \in \mathcal{Y}$  where  $\mathcal{Y} := \{-1, +1\}$ , and prediction rule  $h_w : \mathbb{R}^d \to \{-1, +1\}$  with

$$(1) h_w(x) = \operatorname{sign}\left(w^{\top}x\right)$$

$$= \operatorname{sign}\left(\sum_{j=1}^{d} w_j x_j\right)$$

$$\operatorname{sign}(\xi) = \begin{cases} -1, & \text{if } \xi < 0 \\ +1, & \text{if } \xi > 0 \end{cases}$$

 $\pm 1$  means either -1 or +1,  $\mathbb{R}_{+}:=(0,+\infty)$ , and  $\left\|x\right\|_{2}:=\sqrt{\sum_{\forall j}\left(x_{j}\right)^{2}}$  for the Euclidean distance.

<sup>&</sup>lt;sup>1</sup>We use standard notation

Let the hypothesis class is

(3) 
$$\mathcal{H} = \left\{ x \to w^{\top} x : \forall w \in \mathbb{R}^d \right\}$$

In other words, the hypothesis  $h_w \in \mathcal{H}$  is parametrized by  $w \in \mathbb{R}^d$ , it receives an input vector  $x \in \mathcal{X} := \mathbb{R}^d$  and it returns the label  $y = \text{sign}(w^\top x) \in \mathcal{Y} := \{\pm 1\}$  where

$$\operatorname{sign}(\xi) = \begin{cases} -1, & \text{if } \xi < 0 \\ +1, & \text{if } \xi > 0 \end{cases}$$

Consider a loss function  $\ell : \mathbb{R}^d \to \mathbb{R}_+$  with

(4) 
$$\ell(w, z = (x, y)) = \max(0, 1 - yw^{\top}x) + \lambda ||w||_{2}^{2}$$

for some given value  $\lambda > 0$ .

Assume there is available a dataset of examples  $S_n = \{z_i = (x_i, y_i); i = 1, ..., n\}$  of size n. Do the following:

(1) Show that the function  $f: \mathbb{R} \to \mathbb{R}_+$  with  $f(x) = \max(0, 1 - x)$  is convex in  $\mathbb{R}$ ; and show that the loss (4) is convex.

**Hint:**: You may use Proposition ?? from Handout ??: Elements of convex learning problems.

(2) Show that the loss  $\ell(w, z)$  for  $\lambda = 0$  (4) is L-Lipschitz (with respect to w) when  $x \in \mathcal{X}$  where  $\mathcal{X} := \{x \in \mathbb{R}^d : ||x||_2 \leq L\}$ .

**Hint:** You may use the definition of Lipschitz function. Without loss of generality, you can consider any  $w_1 \in \mathbb{R}^d$  and  $w_2 \in \mathbb{R}^d$  such that  $1 - yw_2^\top x \le 1 - yw_1^\top x$ , and then take cases  $1 - yw_2^\top x > \text{or} < 0$  and  $1 - yw_1^\top x > \text{or} < 0$  to deal with the max.

(3) Construct the set of sub-gradients  $\partial f(x)$  for  $x \in \mathbb{R}$  of the function  $f: \mathbb{R} \to \mathbb{R}_+$  with  $f(x) = \max(0, 1-x)$ . Show that the vector v with

$$v = \begin{cases} 2\lambda w, & yw^{\top}x > 1\\ 2\lambda w, & yw^{\top}x = 1\\ -yx + 2\lambda w, & yw^{\top}x < 1 \end{cases}$$

is  $v \in \partial_w \ell(w, z = (x, y))$ , aka a sub-gradient of  $\ell(w, z = (x, y))$  at w, for any  $w \in \mathbb{R}^d$ .

(4) Write down the algorithm of online AdaGrad (Adaptive Stochastic Gradient Descent) with learning rate  $\eta_t > 0$ , batch size m, and termination criterion  $t > T_{\text{max}}$  for some  $T_{\text{max}} > 0$  in order to discover  $w^*$  such as

(5) 
$$w^* = \arg\min_{\forall w: h_w \in \mathcal{H}} \left( \mathbb{E}_{z \sim g} \left( \ell \left( w, z = (x, y) \right) \right) \right)$$

The formulas in your algorithm should be implemented for the above learning problem and tailored to 1, 3, and 4.

(5) Use the R code given below in order to generate the dataset of observed examples  $S_n = \{z_i = (x_i, y_i)\}_{i=1}^n$  that contains  $n = 10^6$  examples with inputs x of dimension d = 2. Consider  $\lambda = 0$ . Use a seed  $w^{(0)} = (0, 0)^{\top}$ .

- (a) By using appropriate values for m,  $\eta_t$  and  $T_{\text{max}}$ , code in R the algorithm you designed in part 4, and run it.
- (b) Plot the trace plots for each of the dimensions of the generated chain  $\{w^{(t)}\}$  against the iteration t.
- (c) Report the value of the output  $w_{\text{adaGrad}}^*$  (any type) of the algorithm as the solution to (5).
- (d) To which cluster y (i.e., -1 or 1)  $x_{\text{new}} = (1,0)^{\top}$  belongs?

```
\# R code. Run it before you run anything else
#
data_generating_model <- function(n,w) {</pre>
z <- rep( NaN, times=n*3 )
z \leftarrow matrix(z, nrow = n, ncol = 3)
z[,1] \leftarrow rep(1,times=n)
z[,2] \leftarrow runif(n, min = -10, max = 10)
p \leftarrow w[1]*z[,1] + w[2]*z[,2] p \leftarrow exp(p) / (1+exp(p))
z[,3] \leftarrow rbinom(n, size = 1, prob = p)
ind <-(z[,3]==0)
z[ind,3] < -1
x < z[,1:2]
y < -z[,3]
return(list(z=z, x=x, y=y))
n_{obs} < 1000000
w_{true} <- c(-3,4)
set.seed(2023)
out <- data_generating_model(n = n_obs, w = w_true)</pre>
set.seed(0)
z_{obs} \leftarrow out$z #z=(x,y)
x \leftarrow \text{out}
y <- out$y
#z_obs2=z_obs
\#z_obs2[z_obs[,3]==-1,3]=0
#w_true <- as.numeric(glm(z_obs2[,3]~ 1+ z_obs2[,2],family = "binomial"</pre>
)$coefficients)
```