

Exercise sheet

Lecturer/Author: Georgios P. Karagiannis

georgios.karagiannis@durham.ac.uk

Part 1. Convex learning problems

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

Solution. Let $u, v \in \mathbb{R}^d$ and $a \in [0, 1]$. It is

$$\begin{aligned}
 f(\alpha u + (1 - \alpha)v) &= g(\langle \alpha u + (1 - \alpha)v, x \rangle + y) \\
 &= g(\langle \alpha u, x \rangle + \langle (1 - \alpha)v, x \rangle + y) \\
 &= g(\alpha \langle u, x \rangle + y + (1 - \alpha) \langle v, x \rangle + y) & y = \alpha y + (1 - \alpha)y \\
 &\leq \alpha g(\langle u, x \rangle + y) + (1 - \alpha) g(\langle v, x \rangle + y) & (g \text{ is convex}) \\
 &= \alpha f(u) + (1 - \alpha) f(v)
 \end{aligned}$$

Exercise 2. (★) Let functions g_1 be ρ_1 -Lipschitz and g_2 be ρ_2 -Lipschitz. Then, show that, f with $f(x) = g_1(g_2(x))$ is $\rho_1\rho_2$ -Lipschitz.

Solution.

$$\begin{aligned}
 |f(w_1) - f(w_2)| &= |g_1(g_2(w_1)) - g_1(g_2(w_2))| \\
 &\leq \rho_1 |g_2(w_1) - g_2(w_2)| \\
 &\leq \rho_1 \rho_2 |w_1 - w_2|
 \end{aligned}$$

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

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

$$f(v) = g(\langle w, x \rangle + y)$$

$$\leq g(\langle w, x \rangle + y) + g'(\langle w, x \rangle + y) \langle v - w, x \rangle + \frac{\beta}{2} (\langle v - w, x \rangle)^2 \quad (g \text{ is smooth})$$

$$\leq g(\langle w, x \rangle + y) + g'(\langle w, x \rangle + y) \langle v - w, x \rangle + \frac{\beta}{2} (\|v - w\| \|x\|)^2 \quad (\text{Cauchy-Schwarz inequality})$$

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

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

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

Solution. \Rightarrow Let $f : S \rightarrow \mathbb{R}$ be ρ -Lipschitz over convex set S , $w \in S$ and $v \in \partial f(w)$.

- Since S is open we get that there exist $\epsilon > 0$ such as $u := w + \epsilon \frac{v}{\|v\|}$ where $u \in S$. So $\langle u - w, v \rangle = \epsilon \|v\|$ and $\|u - w\| = \epsilon$.
- From the subgradient definition we get

$$f(u) - f(w) \geq \langle u - w, v \rangle = \epsilon \|v\|$$

- From the Lipschitzness of $f(\cdot)$ we get

$$f(u) - f(w) \leq \rho \|u - w\| = \rho \epsilon$$

Therefore $\|v\| \leq \rho$.

\Leftarrow It is for all $w \in S$ and $v \in \partial f(w)$ it is $\|v\| \leq \rho$.

- For any $u \in S$, it is

$$\begin{aligned} f(w) - f(u) &\leq \langle v, w - u \rangle && (\text{because } v \in \partial f(w)) \\ (0.1) \quad &\leq \|v\| \|w - u\| && \text{by Cauchy-Schwarz inequality} \\ &\leq \rho \|w - u\| && \text{because } \|v\| \leq \rho \end{aligned}$$

- Similarly it results $u, w \in S$

$$f(w) - f(u) \leq \langle v, u - w \rangle \|v\| \leq \|v\| \|u - w\| \leq \rho \|u - w\|$$

from (0.1) because w, u can be swaped in (0.1) as they both are any values in S .

Exercise 5. (★) Let $g_1(w), \dots, g_r(w)$ be r convex functions, and let $f(\cdot) = \max_{j \in \{1, \dots, r\}} (g_j(\cdot))$. 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 .

Solution. Since g_k is convex, for all u

$$g_k(u) \geq g_k(w) + \langle u - w, \nabla g_k(w) \rangle$$

However $f(u) = \max_{\forall j} (g_j(u)) \geq g_k(u)$ for any j , and $f(w) = g_k(w)$ at w . Then

$$\begin{aligned} f(u) &\geq g_k(u) \\ &\geq g_k(w) + \langle u - w, \nabla g_k(w) \rangle \\ &= f(w) + \langle u - w, \nabla g_k(w) \rangle \end{aligned}$$

Then by the definition of the sub-gradient $\nabla g_k(w) \in \partial f(w)$

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 Lipschitzness.

Solution. According to the definitions given in the lecture:

- Convex-Lipschitz-Bounded Learning Problem $(\mathcal{H}, \mathcal{Z}, \ell)$ with parameters ρ , and B , is called the learning problem whose the hypothesis class \mathcal{H} is a convex set, for all $w \in \mathcal{H}$ it is $\|w\| \leq B$, and the loss function $\ell(\cdot, z)$ is convex and ρ -Lipschitz function for all $z \in \mathcal{Z}$.

I have:

Convexity: The function $g : \mathbb{R} \rightarrow \mathbb{R}$, defined by $g(a) = a^2$ is convex. Eg. $\frac{d^2}{da^2}g(a) = 2 \geq 0$ is non-negative. The convexity of $\ell(w, z = (x, y))$ for all z follows as a composition of g with a linear function.

Lipschitzness: The function $g(a) = a^2$ is 1-Lipschitz since It is

$$|g(a_2) - g(a_1)| = |a_2^2 - a_1^2| = |(a_2 + a_1)(a_2 - a_1)| \leq 2\rho(a_2 - a_1) = 2\rho|a_2 - a_1|$$

Hence because $|x| \leq \rho$, $g(a)$ is $2\rho^2$ -Lipschitz as a composition.

Boundness: The norm of each hypothesis w is bounded by ρ according to the assumptions.

Therefore,

- (1) the learning problem under consideration is a Convex-Lipschitz-Bounded learning problem.
- (2) the parameter of Lipschitzness is $2\rho^2$.

Exercise 7. (★) If f is λ -strongly convex and u is a minimizer of f then for any w

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

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

Solution.

Exercise 8. (★) Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex and β -smooth function.

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

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

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

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

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

$$\|\nabla f(w)\| \leq \sqrt{2\beta f(w)}$$

Solution.

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

$$w^* = \arg \min_w \{f(w)\}$$

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

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

Hint: Use the definition of λ -strongly convex function, properly rearrange it, and ...

Solution.

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

Solution. We just need to check that for all w, u , and $\alpha \in (0, 1)$ we have

$$\begin{aligned} J(\alpha w + (1 - \alpha)u; \lambda) &\leq \alpha J(w; \lambda) + (1 - \alpha) J(u; \lambda) - \frac{2\lambda}{2} \alpha (1 - \alpha) \|w - u\|^2 \iff \\ \|\alpha w + (1 - \alpha)u\|_2^2 &\leq \alpha \|w\|_2^2 + (1 - \alpha) \|u\|_2^2 - \alpha (1 - \alpha) \|w - u\|_2^2 \iff 0 \leq 0 \end{aligned}$$

Exercise 11. (★★) Consider a learning problem $(\mathcal{H}, \mathcal{Z}, \ell)$ with $\mathcal{H} \subset \mathbb{R}^d$, $d > 0$, and loss function $\ell : \mathcal{H} \times \mathcal{Z} \rightarrow \mathbb{R}_+$ which is convex, β -smooth and non-negative. Let \mathfrak{A} be a learning algorithm with output $\mathfrak{A}(\mathcal{S})$ trained against training dataset $\mathcal{S} = \{z_1, \dots, z_m\}$ of IID samples $z_1, \dots, z_m \sim g$

where g is a data generating distribution. In particular, consider that $\mathfrak{A}(\mathcal{S})$ is the Regularized Loss Minimization learning rule that outputs a hypothesis in

$$\min_w \left\{ \hat{R}_{\mathcal{S}}(w) + \lambda \|w\|_2^2 \right\}$$

for $\lambda \geq \frac{2\beta}{m}$ where $\hat{R}_{\mathcal{S}}(w) = \frac{1}{m} \sum_{i=1}^m \ell(w, z_i)$ for all $w \in \mathcal{H}$.

(1) Prove that

$$\mathbb{E}_{\mathcal{S} \sim g} \left(\hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) \right) \leq R_g(w) + \lambda \|w\|_2^2$$

for all $w \in \mathcal{H}$. $R_g(\cdot)$ denotes the risk function under the real data generating distribution g .

(2) Prove that

$$\mathbb{E}_{\mathcal{S} \sim g} \left(R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) \right) \leq \frac{48\beta}{\lambda m} \mathbb{E}_{\mathcal{S} \sim g} \left(\hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) \right).$$

Hint:: If needed you can use the following:

Let $\mathcal{S}^{(i)} = \{z_1, \dots, z_{i-1}, z', z_{i+1}, \dots, z_m\}$ be a set resulting from \mathcal{S} by replacing its i -th element z_i with an independently drawn $z' \sim g$. Then

$$24\beta \ell(\mathfrak{A}(\mathcal{S}), z_i) + \lambda m \ell(\mathfrak{A}(\mathcal{S}), z_i) + 24\beta \ell(\mathfrak{A}(\mathcal{S}^{(i)}), z') - \lambda m \ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) \geq 0$$

(3) Show that the learning algorithm \mathfrak{A} is on-average-replace-one-stable with rate ε . Specify that rate ε as a function of β , λ , m and possibly any other user specified constants if needed. Explain how the shrinkage parameter λ , the training dataset size m , and the smoothness parameter β affect the stability of the learning algorithm \mathfrak{A} .

(4) Show that the expected risk is bounded as follows

$$\mathbb{E}_{\mathcal{S} \sim g} (R_g(\mathfrak{A}(\mathcal{S}))) \leq \left(1 + \frac{48\beta}{\lambda m} \right) \left(R_g(w) + \lambda \|w\|_2^2 \right)$$

for all $w \in \mathcal{H}$.

Solution.

(1) We have

$$\begin{aligned} \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) &\leq \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) + \lambda \|\mathfrak{A}(\mathcal{S})\|_2^2 \\ &\leq \hat{R}_{\mathcal{S}}(w) + \lambda \|w\|_2^2; \quad \forall w \in \mathcal{W} \end{aligned}$$

and by taking expectations w.r.t. \mathcal{S} , it is

$$(0.2) \quad \mathbb{E}_{\mathcal{S} \sim g} \left(\hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) \right) \leq R_g(w) + \lambda \|w\|_2^2; \quad \forall w \in \mathcal{W}$$

because $\mathbb{E}_{\mathcal{S} \sim g} \left(\hat{R}_{\mathcal{S}}(\cdot) \right) = \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\mathcal{S} \sim g} (\ell(\cdot, z_i)) = R_g(\cdot)$.

(2) From a well known theorem to us, it is

$$\mathbb{E}_{\mathcal{S} \sim g} \left(R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})) \right) = \mathbb{E}_{\mathcal{S}, z', i} \left(\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i) \right)$$

Now I ll gonna work on the second term as this is what I ve given in the hint...

It is

$$24\beta\ell(\mathfrak{A}(\mathcal{S}), z_i) + \lambda m\ell(\mathfrak{A}(\mathcal{S}), z_i) + 24\beta\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z') - \lambda m\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) \geq 0 \Leftrightarrow$$

$$\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i) \leq \frac{24\beta}{\lambda m} (\ell(\mathfrak{A}(\mathcal{S}), z_i) + \ell(\mathfrak{A}(\mathcal{S}^{(i)}), z'))$$

Taking expectations

$$\mathbb{E}_{\mathcal{S}, z', i} (\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i)) \leq \frac{24\beta}{\lambda m} \mathbb{E}_{\mathcal{S}, z', i} (\ell(\mathfrak{A}(\mathcal{S}), z_i) + \ell(\mathfrak{A}(\mathcal{S}^{(i)}), z'))$$

Due to the sampling it is

$$\mathbb{E}_{\mathcal{S}} (\ell(\mathfrak{A}(\mathcal{S}), z_i)) = \mathbb{E}_{\mathcal{S}, z', i} (\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z')) = \mathbb{E}_{\mathcal{S} \sim g} (\hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})))$$

So I get

$$\mathbb{E}_{\mathcal{S}, z', i} (\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i)) \leq \frac{48\beta}{\lambda m} \mathbb{E}_{\mathcal{S} \sim g} (\hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S}))).$$

So I get

$$\mathbb{E}_{\mathcal{S} \sim g} (R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S}))) \leq \frac{48\beta}{\lambda m} \mathbb{E}_{\mathcal{S} \sim g} (\hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S}))).$$

(3) Okay, let's say, I did not do the previous part. I see it is

$$\mathbb{E}_{\mathcal{S} \sim g} (R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S}))) \leq \frac{48\beta}{\lambda m} \mathbb{E}_{\mathcal{S} \sim g} (\hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S}))).$$

From a well known theorem to us, it is

$$\mathbb{E}_{\mathcal{S} \sim g} (R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S}))) = \mathbb{E}_{\mathcal{S}, z', i} (\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i))$$

So

$$\mathbb{E}_{\mathcal{S}, z', i} (\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i)) \leq \frac{48\beta}{\lambda m} \mathbb{E}_{\mathcal{S} \sim g} (\hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})))$$

but the expectation depends on m ... so

$$\begin{aligned} \mathbb{E}_{\mathcal{S}, z', i} (\ell(\mathfrak{A}(\mathcal{S}^{(i)}), z_i) - \ell(\mathfrak{A}(\mathcal{S}), z_i)) &\leq \frac{48\beta}{\lambda m} \mathbb{E}_{\mathcal{S} \sim g} (\hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S}))) \\ &\leq \frac{48\beta}{\lambda m} \mathbb{E}_{\mathcal{S} \sim g} (\hat{R}_{\mathcal{S}}(0)) \\ &\leq \frac{48\beta}{\lambda m} \mathbb{E}_{\mathcal{S} \sim g} (\max \hat{R}_{\mathcal{S}}(0)) \\ &\leq \frac{48\beta}{\lambda m} \mathbb{E}_{\mathcal{S} \sim g} \left(\frac{1}{m} \sum_{i=1}^n \underbrace{\max \ell(0, z_i)}_{=C} \right) \\ &\leq \frac{48\beta}{\lambda m} C \end{aligned}$$

so it is the learning algorithm \mathfrak{A} is on-average-replace-one-stable with rate

$$\varepsilon = \frac{48\beta}{\lambda m} C$$

and $C = \max \ell(0, z_i)$...or whatever constant they pick.

Larger training sample size m , and larger regularization parameter λ (eg more parsimonious model) lead to a more stable learning algorithm. Smaller smoothness parameter (the gradient changes less wrt the argument) leads to more stable learning algorithm.

(4) We use the decomposition discussed in the lectures,

$$\begin{aligned} \mathbb{E}_{\mathcal{S} \sim g} (R_g(\mathfrak{A}(\mathcal{S}))) &= \underbrace{\mathbb{E}_{\mathcal{S} \sim g} (\hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})))}_{\text{}} + \underbrace{\mathbb{E}_{\mathcal{S} \sim g} (R_g(\mathfrak{A}(\mathcal{S})) - \hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S})))}_{\text{}} \\ &\leq \mathbb{E}_{\mathcal{S} \sim g} (\hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S}))) + \frac{48\beta}{\lambda m} \mathbb{E}_{\mathcal{S} \sim g} (\hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S}))) \\ &= \left(1 + \frac{48\beta}{\lambda m}\right) \mathbb{E}_{\mathcal{S} \sim g} (\hat{R}_{\mathcal{S}}(\mathfrak{A}(\mathcal{S}))) \\ &\leq \left(1 + \frac{48\beta}{\lambda m}\right) (R_g(w) + \lambda \|w\|_2^2), \quad \forall w \in \mathcal{H} \end{aligned}$$

Part 2. Stochastic learning

Exercise 12. (★) Assume a Bayesian model

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

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

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

by using SGD with update

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

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

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

Solution. It is

$$\begin{aligned}
\mathbb{E}_{\mathcal{J}^{(t)} \sim \text{SRS}} \left(\frac{n}{m} \sum_{j \in \mathcal{J}^{(t)}} \nabla_w \log \left(f \left(z_j | w^{(t)} \right) \right) \right) &= \frac{n}{m} \sum_{j \in \mathcal{J}^{(t)}} \mathbb{E}_{\mathcal{J}^{(t)} \sim \text{SRS}} \left(\nabla_w \log \left(f \left(z_j | w^{(t)} \right) \right) \right) \\
&= \frac{n}{m} \sum_{j \in \mathcal{J}^{(t)}} \mathbb{E}_{\mathcal{J}^{(t)} \sim \text{SRS}} \left(\nabla_w \log \left(f \left(z_j | w^{(t)} \right) \right) \right) \\
&= \frac{n}{m} \sum_{j \in \mathcal{J}^{(t)}} \frac{1}{n} \sum_{i=1}^n \nabla_w \log \left(f \left(z_i | w^{(t)} \right) \right) \\
&= \sum_{i=1}^n \nabla_w \log \left(f \left(z_i | w^{(t)} \right) \right)
\end{aligned}$$

It is $\mathbb{E}_{\mathcal{J}^{(t)} \sim \text{SRS}} \left(\nabla_w \log \left(f \left(z_j | w^{(t)} \right) \right) \right) = \frac{1}{n} \sum_{i=1}^n \nabla_w \log \left(f \left(z_i | w^{(t)} \right) \right)$ because the expectation is under the probability I get randomly an integer and for the j th on the probability is $1/n$ due to the random scheme. Also $|\mathcal{J}^{(t)}| = m$.

Exercise 13. (★) Let $\{v_t; t = 1, \dots, T\}$ be a sequence of vectors with $v_t \in \mathbb{R}^d$ and $d \in \mathbb{N} - \{0\}$. Consider an algorithm producing $\{w^{(t)}; t = 1, 2, 3, \dots\}$ with

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

$w_t \in \mathbb{R}^d$ and $d \in \mathbb{N} - \{0\}$. Show that

(1) it is

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

Hint:: Recall that

$$\|x + y\|_2^2 = \|x\|_2^2 + \|y\|_2^2 + 2 \langle x, y \rangle, \quad \forall x, y \in \mathbb{R}^d, d \in \mathbb{N} - \{0\}$$

(2) it is

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

(3) (continue) it is

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

Solution.

(1) It is

$$\begin{aligned}\langle w^{(t)} - w^*, v_t \rangle &= \frac{1}{\eta} \langle w^{(t)} - w^*, \eta v_t \rangle \\ &= \frac{1}{\eta} \left(- \langle w^{(t)} - w^*, -\eta v_t \rangle \right)\end{aligned}$$

Then by using the Hint as

$$\langle x, y \rangle = \frac{1}{2} \left(\|x + y\|_2^2 - \|x\|_2^2 - \|y\|_2^2 \right)$$

for $x = w^{(t)} - w^* \in \mathbb{R}^d$ and $y = -\eta v_t \in \mathbb{R}^d$, I get

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

(2) So

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

(3) So

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

Exercise 14. (★) Let $\{v_t; t = 1, \dots, T\}$ be a sequence of vectors. Consider an algorithm producing $\{w^{(t)}; t = 1, 2, 3, \dots\}$ with

$$\begin{aligned}w^{(1)} &= 0 \\ w^{(t+\frac{1}{2})} &= w^{(t)} - \eta v_t \\ w^{(t+1)} &= \arg \min_{w \in \mathcal{H}} \left(\|w - w^{(t+\frac{1}{2})}\| \right)\end{aligned}$$

for $t = 1, \dots, T$.

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.

$$v = \arg \min_{x \in \mathcal{H}} \|x - w\|^2$$

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

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

Show that

(1) it is

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

(2) it is

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

(3) (continue) it is

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

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.

Solution.

(1) It is

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

because from the Projection Lemma

$$\|w^{(t+1)} - w^*\|^2 \leq \|w^{(t+\frac{1}{2})} - w^*\|^2$$

(2) So

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

(3) So

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

Exercise 15. (★) ¹Consider the binary classification problem with inputs $x \in \mathcal{X}$ where $\mathcal{X} := \{x \in \mathbb{R}^d : \|x\|_2 \leq L\}$ for some given value $L > 0$, target $y \in \mathcal{Y}$ where $\mathcal{Y} := \{-1, +1\}$, and prediction rule $h_w : \mathbb{R}^d \rightarrow \{-1, +1\}$ with

$$(0.3) \quad h_w(x) = \text{sign}(w^\top x)$$

$$(0.4) \quad = \text{sign}\left(\sum_{j=1}^d w_j x_j\right)$$

Let the hypothesis class is

$$(0.5) \quad \mathcal{H} = \{x \rightarrow w^\top x : \forall w \in \mathbb{R}^d\}$$

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

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

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

$$(0.6) \quad \ell(w, z = (x, y)) = \max(0, 1 - yw^\top x) + \lambda \|w\|_2^2$$

¹We use standard notation

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

± 1 means either -1 or $+1$, $\mathbb{R}_+ := (0, +\infty)$, and $\|x\|_2 := \sqrt{\sum_{j=1}^d (x_j)^2}$ for the Euclidean distance.

for some given value $\lambda > 0$.

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

Do the following:

- (1) Show that the function $f : \mathbb{R} \rightarrow \mathbb{R}_+$ with $f(x) = \max(0, 1 - x)$ is convex in \mathbb{R} ; and show that the loss (0.6) is convex.

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

- (2) Show that the loss $\ell(w, z)$ for $\lambda = 0$ (0.6) 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 $1 - yw_2^\top x \leq 1 - yw_1^\top x$, 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 $f : \mathbb{R} \rightarrow \mathbb{R}_+$ 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 $t > T_{\max}$ for some $T_{\max} > 0$ in order to discover w^* such as

$$(0.7) \quad w^* = \arg \min_{w: h_w \in \mathcal{H}} (\mathbb{E}_{z \sim g} (\ell(w, z = (x, y))))$$

The formulas in your algorithm should be implemented for the above learning problem and tailored to 0.3, 0.5, and 0.6.

- (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 , η_t and T_{\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_{adaGrad}^* (any type) of the algorithm as the solution to (0.7).
 - (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) {
  z <- rep( NaN, times=n*3 )
  z <- matrix(z, nrow = n, ncol = 3)
  z[,1] <- rep(1,times=n)
  z[,2] <- runif(n, min = -10, max = 10)
  p <- w[1]*z[,1] + w[2]*z[,2] p <- exp(p) / (1+exp(p))
  z[,3] <- 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)
set.seed(0)
z_obs <- out$z #z=(x,y)
x <- out$x
y <- out$y
#z_obs2=z_obs
#z_obs2[z_obs[,3]==-1,3]=0
#w_true <- as.numeric(glm(z_obs2[,3]~ 1+ z_obs2[,2],family = "binomial"
)$coefficients)

```

Solution.

Part 3. Support Vector Machines

Part 4. The kernel trick

Part 5. Multi-class classification

Part 6. Artificial Neural Networks

Part 7. Gaussian process regression

Part 8. Revision