

## 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  $\rho_1$ -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.

**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  $\beta$ -smooth function. Then show that  $f$  is a  $(\beta \|x\|^2)$ -smooth.

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

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

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

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

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

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

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

**Solution.**  $\Rightarrow$  Let  $f : S \rightarrow \mathbb{R}$  be  $\rho$ -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

$$\begin{aligned} f(w) - f(u) &\leq \langle v, w - u \rangle && (\text{because } v \in \partial f(w)) \\ (1) \quad &\leq \|v\| \|w - u\| && \text{by Cauchy-Schwarz inequality} \\ &\leq \rho \|w - u\| && \text{because } \|v\| \leq \rho \end{aligned}$$

- 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 (1) because  $w, u$  can be swapped in (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  $\rho$ , 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  $\rho$ -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  $\rho$  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  $\lambda$ -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.**

---

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

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

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

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

(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.**

(1) If  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is  $\beta$ -smooth then it is

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

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

If it is convex then it is

$$f(v) \geq f(w) + \langle \nabla f(w), v - w \rangle$$

$$f(v) - f(w) \geq \langle \nabla f(w), v - w \rangle$$

Together these conditions imply upper and lower bounds

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

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

$$f(v) \leq f(w) + \langle \nabla f(w), v - w \rangle + \frac{\beta}{2} \|v - w\|_2^2 \quad (\text{due to smoothness})$$

$$\iff f(w) - f(v) \leq f(w) - f(v)$$

$$\iff \langle \nabla f(w), v - w \rangle + \frac{\beta}{2} \|v - w\|_2^2 \leq f(w) - f(v)$$

$$\iff \left\langle \nabla f(w), \frac{1}{\beta} \nabla f(w) \right\rangle + \frac{\beta}{2} \left\| \frac{1}{\beta} \nabla f(w) \right\|_2^2 \leq f(w) - f(v)$$

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

$$\|\nabla f(w)\|^2 \leq 2\beta (f(w) - f(v))$$

as  $f(\cdot) \geq 0$

$$\|\nabla f(w)\|^2 \leq 2\beta f(w)$$

(3) From part 2, this is obvious because  $f(x) > 0$  for all  $x \in \mathbb{R}^d$ , as

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

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

**Exercise 9.** (★) Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be a  $\lambda$ -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  $\lambda$ -strongly convex function, properly rearrange it, and ...

**Solution.** We use the definition of  $\lambda$ -strongly convex function; i.e. for all  $w, u$ , and  $\alpha \in (0, 1)$  we have

$$\begin{aligned} f(\alpha w + (1 - \alpha)u) &\leq \alpha f(w) + (1 - \alpha)f(u) - \frac{\lambda}{2} \alpha(1 - \alpha) \|w - u\|^2 \Leftrightarrow \\ \frac{f(\alpha w + (1 - \alpha)u) - f(u)}{\alpha} &\leq f(w) - f(u) - \frac{\lambda}{2} (1 - \alpha) \|w - u\|^2 \end{aligned}$$

For  $u = w^*$  it is

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

When  $a \rightarrow 0$

$$\frac{\lambda}{2} \alpha(1 - \alpha) \|w - w^*\|^2 \rightarrow 0$$

I know that  $w^*$  is the minimizer of  $f$ . So 0 is the minimizer of  $g$  with  $g(a) = f(\alpha w + (1 - \alpha)w^*)$  hence when  $a \rightarrow 0$

$$\frac{f(\alpha w + (1 - \alpha)w^*) - f(w^*)}{\alpha} \rightarrow \left. \frac{d}{d\alpha} g(\alpha) \right|_{\alpha=0}$$

So

$$0 \leq f(w) + f(w^*) - \frac{\lambda}{2} \|w - w^*\|^2$$

which concludes the proof.

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

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

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

---

**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,  $\beta$ -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  $\varepsilon$ . Specify that rate  $\varepsilon$  as a function of  $\beta$ ,  $\lambda$ ,  $m$  and possibly any other user specified constants if needed. Explain how the shrinkage parameter  $\lambda$ , the training dataset size  $m$ , and the smoothness parameter  $\beta$  affect the stability of the learning algorithm  $\mathfrak{A}$ .

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

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

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

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

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

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

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

It is

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

Taking expectations

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

Due to the sampling it is

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

So I get

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

So I get

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

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

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

From a well known theorem to us, it is

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

So

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

but the expectation depends on  $m$ ... so

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

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  $\lambda$  (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.** (★) 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 13.** (★) 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 17 from “Handout 4: 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 4 in “Handout 4: 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 3 in “Handout 5: 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 14.** (★) <sup>1</sup>Consider the binary classification problem with inputs  $x \in \mathcal{X}$  where  $\mathcal{X} := \{x \in \mathbb{R}^d : \|x\|_2 \leq L\}$  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

$$(3) \quad h_w(x) = \text{sign}(w^\top x)$$

$$(4) \quad = \text{sign}\left(\sum_{j=1}^d w_j x_j\right)$$

Let the hypothesis class is

$$(5) \quad \mathcal{H} = \{x \rightarrow w^\top x : \forall w \in \mathbb{R}^d\}$$

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

$$(6) \quad \ell(w, z = (x, y)) = \max(0, 1 - yw^\top x) + \lambda \|w\|_2^2$$

for some given value  $\lambda > 0$ .

Assume there is available a dataset of examples  $S_n = \{z_i = (x_i, y_i) ; i = 1, \dots, n\}$  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 (6) is convex.

**Hint::** You may use Note 11 from Lecture notes 2: Elements of convex learning problems.

- (2) Show that the loss  $\ell(w, z)$  for  $\lambda = 0$  (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.

---

<sup>1</sup>We use standard notation

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

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

- (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

$$(7) \quad w^* = \arg \min_{w: h_w \in \mathcal{H}} (\mathbb{E}_{z \sim g} (\ell(w, z = (x, y))))$$

The formulas in your algorithm should be implemented for the above learning problem and tailored to 3, 5, and 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  $m$ ,  $\eta_t$  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  $w_{\text{adaGrad}}^*$  (any type) of the algorithm as the solution to (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.**

- (1)  $f_1(x) = 0$  is convex,  $f_2(x) = 1 - x$  is convex, hence from the example in the lecture notes,  $f(x) = \max(f_1(x), f_2(x))$  is convex as well. Regarding the loss function, we just have  $f_2(w) = 1 - yx^\top w$  which is convex as a composition due to linearity.
- (2) Given a fixed example  $(x, y) \in \{x \in \mathbb{R}^d : \|x'\|_2 \leq R\} \times \{-1, 1\}$ .  
 Assume  $w_1, w_2 \in \mathbb{R}^d$ . Let  $\ell_i = \max\{0, 1 - yx^\top w_i\}$ , for  $i = 1, 2$ . It suffices to show that  $|\ell_1 - \ell_2|_2 \leq R|w_1 - w_2|_2$ . I take cases  
**Case-1:** Assume  $yx^\top w_1 \geq 1$  and  $yx^\top w_2 \geq 1$  then  $|\ell_1 - \ell_2|_2 = 0 \leq R|w_1 - w_2|_2$

**Case-2:** Assume that at least one of  $yx^\top w_1 < 1$  or  $yx^\top w_2 < 1$  but not both is true.  
Assume without loss of generality that  $1 - yx^\top w_1 < 1 - yx^\top w_2$ . Then

$$\begin{aligned}
|\ell_1 - \ell_2|_2 &= \ell_1 - \ell_2 \\
&= 1 - yx^\top w_1 - \max(0, 1 - yx^\top w_2) \\
&\leq 1 - yx^\top w_1 - (1 - yx^\top w_2) \\
&= yx^\top (w_2 - w_1) \\
&\leq y \|x^\top\|_2 \|w_1 - w_2\|_2 \quad \text{because } a^\top b \leq \|a\| \|b\|
\end{aligned}$$

(3) It is

$$f(x) = \max(0, 1 - x) = \begin{cases} 0 & x > 1 \\ 0 & x = 1 \\ 1 - x & x < 1 \end{cases}$$

- For  $x > 1$ ,  $f$  is differentiable so  $\partial f(x) = \{f'(x)\} = \{0\}$ .
- For  $x < 1$ ,  $f$  is differentiable so  $\partial f(x) = \{f'(x)\} = \{-1\}$ .
- For  $x = 1$ ,  $f$  is not differentiable. By definition I have that  $v$  is subgradient of  $f(x)$  at  $x = 0 \in S$  if

$$\forall u \in \mathbb{R}, \quad f(u) \geq f(x) + \langle u - x, v \rangle$$

So, for  $u \geq 1$ , it is  $0 \geq (u - 1)v \implies v \leq 0$ , and for  $u < 1$  it is  $(1 - u) \geq (u - 1)v \implies v \geq -1$ . Hence the common space is  $v \in [0, 1]$  So  $\partial f(x) = [0, 1]$ . Hence,

$$\partial f(x) = \begin{cases} 0, & x > 1 \\ [-1, 0], & x = 1 \\ -1, & x < 1 \end{cases}$$

Now regarding the loss  $\partial_w \ell(w, z = (x, y))$

- for  $yw^\top x > 1$  it is differentiable so  $\nabla_w \ell(w, z = (x, y)) = \nabla_w (0 + \lambda \sum_{j=1}^d w_j^2) = 2\lambda w$ ;  
as

$$\frac{d}{dw_j} \sum_{j'=1}^d w_{j'}^2 = 2\lambda w_j$$

- for  $yw^\top x < 1$  it is differentiable so  $\nabla_w \ell(w, z = (x, y)) = \nabla_w (1 - yw^\top x + \lambda \sum_{j=1}^d w_j^2) = yx + 2\lambda w$  as

$$\frac{d}{dw_j} (1 - yw^\top x) = \frac{d}{dw_j} \left( 1 - y \sum_{j'=1}^d w_{j'} x_{j'} \right) = -yx_j$$

- for  $yw^\top x = 1$ ,  $v = 0$  satisfies the definition of the sub-gradient

$$\begin{aligned} \forall u, f(u) &\geq \overset{0}{f(w)} + \langle u - w, v \rangle \\ \max(0, 1 - yu^\top x) &\geq 0 + (u - w)^\top 0 \end{aligned}$$

So

$$\begin{aligned} \partial \ell(w, z = (x, y)) &= \partial \left( \max(0, 1 - yw^\top x) + \lambda \|w\|_2^2 \right) \\ &= \partial \left( \max(0, 1 - yw^\top x) \right) + \partial \left( \lambda \|w\|_2^2 \right) \\ &= \partial \left( \max(0, 1 - yw^\top x) \right) + \nabla \left( \lambda \|w\|_2^2 \right) \\ &= 0 + 2\lambda w \end{aligned}$$

but  $\partial \left( \lambda \|w\|_2^2 \right) = \left\{ \nabla \left( \lambda \|w\|_2^2 \right) \right\}$  because  $\lambda \|w\|_2^2$  is differentiable. Hence

$$\partial \ell(w, z = (x, y)) = 0 + 2\lambda w$$

Hence

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

(4)

**Algorithm.** For  $t = 1, 2, 3, \dots$  iterate:

- Get a random sub-sample  $\left\{ \tilde{z}_i^{(t)} = (\tilde{x}_i^{(t)}, \tilde{y}_i^{(t)}) ; i = 1, \dots, m \right\}$  of size  $m$  with or without replacement from the complete data-set  $\mathcal{S}_n$ .
- For  $j = 1, \dots, d$  (index  $j$  indicates the dimension of  $w$ ) compute

$$w_j^{(t+1)} = w_j^{(t)} - \eta_t \frac{1}{\sqrt{[G_t]_{j,j} + \epsilon}} \bar{v}_{t,j}$$

$[G_t]_{j,j} = [G_{t-1}]_{j,j} + (\bar{v}_{t,j})^2$  where  $\bar{v}_t = \frac{1}{m} \sum_{i=1}^m \tilde{v}_{t,i}$  and

$$\tilde{v}_{t,i} = \begin{cases} 2\lambda w^{(t)}, & \tilde{y}_i^{(t)} (w^{(t)})^\top \tilde{x}_i^{(t)} > 1 \\ 2\lambda w^{(t)}, & \tilde{y}_i^{(t)} (w^{(t)})^\top \tilde{x}_i^{(t)} = 1 \\ -\frac{1}{m} \tilde{y}_i^{(t)} \tilde{x}_i^{(t)} + 2\lambda w^{(t)}, & \tilde{y}_i^{(t)} (w^{(t)})^\top \tilde{x}_i^{(t)} < 1 \end{cases}$$

where index  $i$  indicates the sub-sample, and  $\epsilon > 0$  small.

- Terminate if a termination criterion is satisfied

(5)

- The R code can be found in the link [https://raw.githubusercontent.com/georgios-stats/Machine\\_Learning\\_and\\_Neural\\_Networks\\_III\\_Epiphany\\_2025/main/Exercise\\_sheets/supplementary/q6\\_adagrad.R](https://raw.githubusercontent.com/georgios-stats/Machine_Learning_and_Neural_Networks_III_Epiphany_2025/main/Exercise_sheets/supplementary/q6_adagrad.R)
- The figures are presented below



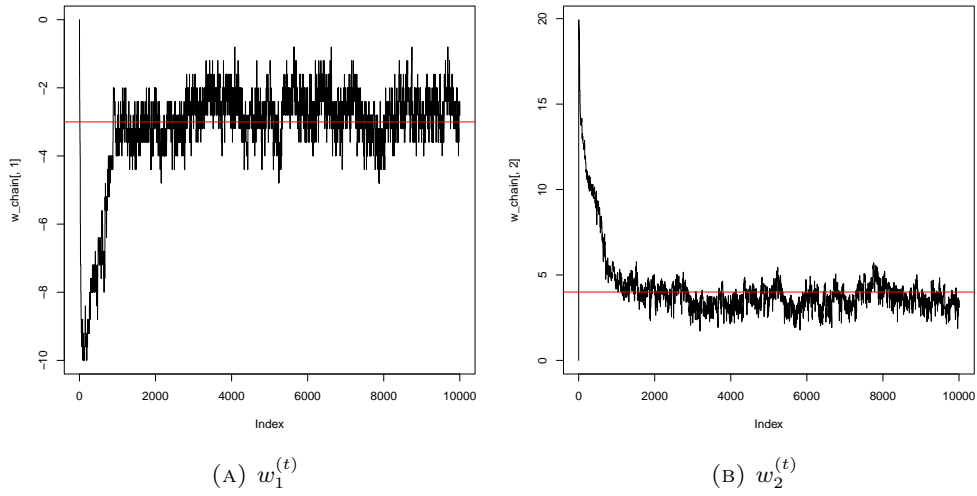


FIGURE 0.1. trace plots

- (c) I found  $w = (-2.674615, 3.205785)$   
 (d) It belongs to  $-1$

**Exercise 15.** (★) 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  $|\mathcal{J}^{(t)}| = m$ .

---

### Part 3. Support Vector Machines

---

**Exercise 16.** (★★) Consider a training data set  $\mathcal{S} = \{z_i = (x_i, y_i)\}_{i=1}^m$ . Consider the Soft-SVM Algorithm that requires the solution of the following quadratic minimization problem (in a slightly modified but equivalent form to what we have discussed)

**Primal problem:**

$$(8) \quad (w^*, b^*, \xi^*) = \arg \min_{(w, b, \xi)} \left( \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^m \xi_i \right)$$

$$(9) \quad \text{subject to: } y_i (\langle w, x_i \rangle + b) \geq 1 - \xi_i, \quad \forall i = 1, \dots, m$$

$$(10) \quad \xi_i \geq 0, \quad \forall i = 1, \dots, m$$

for some user-specified fixed parameter  $C > 0$ .

- (1) Specify the Lagrangian function  $L$  associated to the above primal quadratic minimization problem, where  $\{\alpha_i\}$  are the Lagrange coefficients wrt (9), and  $\{\beta_i\}$  are the Lagrange coefficients wrt (10). Write down any possible restrictions on the Lagrange coefficients.
- (2) Compute the dual Lagrangian function denoted as  $\tilde{L}$  as a function of the Lagrange coefficients and the data points  $\mathcal{S}$ .
- (3) Apply the Karush–Kuhn–Tucker (KKT) conditions to the above problem, and write them down.
- (4) Derive and write down the dual Lagrangian quadratic maximization problem, along with the inequality and equality constraints, where you seek to find  $\{\alpha_i\}$ .
- (5) Justify why the  $i$ -th point  $x_i$  lies on the margin boundary when  $\alpha_i \in (0, C)$  (beware it is  $\alpha_i \neq C$ ), and why the  $i$ -th point  $x_i$  lies inside the margin when  $\alpha_i = C$ .
- (6) Given optimal values  $\{\alpha_i^*\}$  for Lagrangian coefficients  $\{\alpha_i\}$  as they are derived by solving the dual Lagrangian maximization problem in part 4, derive the optimal values  $w^*$  and  $b^*$  for the parameters  $w$  and  $b$  as function of the support vectors. Regarding parameter  $b$  it should be in the derived in the form

$$b^* = \frac{1}{|\mathcal{M}|} \sum_{i \in \mathcal{M}} \left( y_i - \sum_{j \in \mathcal{S}} \alpha_j^* y_j \langle x_j, x_i \rangle \right)$$

where you determine the sets  $\mathcal{M}$  and  $\mathcal{S}$ .

- (7) Report the halfspace predictive rule  $h_{w,b}(x)$  of the above problem as a function of  $\alpha^*$  and  $b^*$ .

**Solution.**

---

**Exercise 17.** (★★★) [This is the Relevance Vector Machine. The Exercise is taken from “Exercise Sheet: Bayesian Statistics” of the module “Bayesian Statistics III/IV (MATH3361/4071)” taught  
Page 19 Created on 2025/02/11 at 16:14:12 by Georgios Karagiannis]

in “Michaelmas term 2021”. The supplementary material in the box was mainly provided for the students who had not been introduced to the SVM ideas or the Kernel trick -so it can be skipped. Also, the supplementary material in the box is presented with a statistical (geostatistical modeling) motivation. The exercise requires basic knowledge of Bayesian statistical inference and in particular the use of Bayes theorem for the computation of the posterior as well as basis probability density calculus. However, the exercise is a useful example of extending the SVM ideas to the Bayesian learning setting./

Regarding the statistical model: Long story, short (supplementary material)

Consider that we are interested in recovering the mapping

$$x \mapsto^{\eta} \eta(x)$$

in the sense that  $y \in \mathbb{R}$  is the response (output quantity) that depends on  $x = (x_1, \dots, x_d) \in \mathcal{X} \subseteq \mathbb{R}^d$  which is the independent variable (input quantity) in a procedure; E.g.:

- $y$ : precipitation in log scale
- $x = (\text{longitude}, \text{latitude})$ : geographical coordinates.

Consider a set of observed data  $\{(y_i, x_i)\}_{i=1}^n$ , which may be contaminated by additive noise of unknown variance; i.e.

$$y_i = \eta(x_i) + \epsilon_i,$$

where  $\epsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$  and  $\sigma^2 > 0$  is unknown. We wish to recover  $\eta(x)$  by using the Tikhonov regularization on the functional space  $\mathcal{H}$  such that

$$(11) \quad \eta = \arg \min_{\forall \tilde{\eta} \in \mathcal{H}} \left\{ \sum_{i=1}^n L(y_i - \tilde{\eta}(x_i)) + \lambda \|\tilde{\eta}\|_{\mathcal{H}}^2 \right\}$$

By assuming that  $\mathcal{H}$  is a Reproducing Kernel Hilbert Space (RKHS), the solution to the above Ridge regularizes loss minimization problem is such that

$$\eta(x) = \beta_0 + \sum_{j=1}^n k(x, x_j) \beta_j = k(x)^\top \beta$$

where  $k(x) = (1, k(x, x_1), \dots, k(x, x_n))^\top$ ,  $k(x, x_j)$  is the reproducing kernel (such as  $k_\phi(x, x_j) = \exp(-\phi \|x - x_j\|^2)$  for some known parameter  $\phi > 0$ ), and  $\beta \in \mathbb{R}^{n+1}$  is an unknown vector.

Consider the following Bayesian model<sup>2</sup>

$$\begin{cases} y|\beta, \sigma^2 & \sim N(K\beta, I\sigma^2) \\ \beta|\lambda & \sim N(0, D^{-1}), \quad D = (\lambda_0, \lambda_1, \dots, \lambda_n) \\ \lambda_i & \stackrel{\text{iid}}{\sim} d\Pi(\lambda_i) \propto \lambda_i^{a-1} \exp(-b\lambda_i) d\lambda_i, \quad \forall i = 1, \dots, n \\ \sigma^2 & \sim d\Pi(\sigma^2) \propto (\sigma^2)^{c-1} \exp(-\frac{1}{\sigma^2}d) d\sigma^2 \\ \beta, \sigma^2 & \text{a priori independent} \end{cases}$$

where  $K$  is a known matrix with size  $n \times (n+1)$  such that

$$K = \begin{bmatrix} 1 & k(x_1, x_1) & \cdots & k(x_1, x_n) \\ \vdots & \vdots & \ddots & \vdots \\ 1 & k(x_n, x_1) & \cdots & k(x_n, x_n) \end{bmatrix}.$$

The quantities  $a > 0$ ,  $b > 0$ ,  $c > 0$ ,  $d > 0$ , and  $\phi > 0$  are considered as fixed.

- (1) When  $b = 0$ , show that a necessary condition for a valid posterior inference is  $a \in (-1/2, 0)$  for any choice of prior for  $\tau$  (i.e. any choice of  $(c, d)$ ).
- (2) Let  $P = K(K^\top K)^{-1}K^\top$ . Show that (2a) and (2b) are sufficient conditions for the Bayesian model to lead to a valid posterior inference
  - (a) if  $a > 0$  and  $b > 0$ , or
  - (b) if  $y^\top (I - P)y + 2d > 0$  and  $c > -\frac{n}{2}$
- (3) Does the the improper Uniform prior on the joint  $\log(\lambda_i)$  and  $\log(\sigma^2)$ , i.e.  $\pi(\log(\lambda_i), \log(\sigma^2)) \propto 1$ , lead to a valid inference?
- (4) Does the Jeffreys' prior  $\pi(\lambda_i) \propto 1/\lambda_i$  lead to a valid inference?

**Hint-1::**

$$(y - K\beta)^\top (y - K\beta) + (\beta - \mu)^\top V^{-1}(\beta - \mu) = (\beta - \mu^*)^\top (V^*)^{-1}(\beta - \mu^*) + S^*;$$

$$S^* = \mu^\top V^{-1}\mu - (\mu^*)^\top (V^*)^{-1}(\mu^*) + y^\top y; \quad V^* = (V^{-1} + K^\top K)^{-1}; \quad \mu^* = V^* (V^{-1}\mu + K^\top y)$$

**Hint-2::** Sherman-Morrison-Woodbury formula:

$$(A + UCV)^{-1} = A^{-1} - A^{-1}U(C^{-1} + VA^{-1}U)^{-1}VA^{-1}$$

**Hint-3::**

$$-\frac{y^\top y}{2\sigma^2} \leq -\frac{y^\top (I\sigma^2 + KD^{-1}K^\top)^{-1}y}{2} \leq -\frac{1}{2\sigma^2}y^\top (I - P)y$$

where  $P = K(K^\top K)^{-1}K$ .

**Hint-4::** It is given that  $\int_{(0,\infty)} \frac{t^{-(a+1)}}{(\xi+t)^{1/2}} dt < \infty$  if and only if  $a \in (-1/2, 0)$ .

**Solution.**

The posterior pdf is given by

$$\pi(\beta, \sigma^2, \lambda|y) = \frac{f(y|\beta, \sigma^2) \pi(\beta, \sigma^2, \lambda)}{f(y)}$$

---

<sup>2</sup>Dixit, A., & Roy, V. (2021). Posterior impropriety of some sparse Bayesian learning models. *Statistics & Probability Letters*, 171, 109039.

and is proper iff  $f(y) < \infty$  where

$$f(y) = \int \left( \int \left( \underbrace{\int f(y|\beta, \sigma^2) \pi(\beta|\lambda) d\beta}_{=f(y|\lambda, \sigma^2)} \right) \pi(\lambda) d\lambda \right) \pi(\sigma^2) d\sigma^2$$

$\underbrace{\hspace{10em}}_{=f(y|\sigma^2)}$

It is

$$\begin{aligned} f(y|\lambda, \sigma^2) &= \int f(y|\beta, \sigma^2) \pi(\beta, \sigma^2) d\beta \\ &= (2\pi)^{-\frac{n+n+1}{2}} \left( \frac{1}{\sigma^2} \right)^{\frac{n}{2}} \det(D)^{\frac{1}{2}} \int \exp \left( -\frac{1}{2\sigma^2} \left( (y - K\beta)^\top (y - K\beta) + \beta^\top (D\sigma^2) \beta \right) \right) d\beta \\ &= (2\pi)^{-\frac{n+n+1}{2}} \left( \frac{1}{\sigma^2} \right)^{\frac{1}{2}} \det(D)^{\frac{1}{2}} \left[ \int \exp \left( -\frac{1}{2\sigma^2} (\beta - \mu^*)^\top V^* (\beta - \mu^*) \right) d\beta \right] \left[ \exp \left( -\frac{1}{2\sigma^2} S^* \right) \right] \end{aligned}$$

Because

$$\begin{aligned} \int \exp \left( -\frac{1}{2\sigma^2} (\beta - \mu^*)^\top V^* (\beta - \mu^*) \right) d\beta &= (2\pi)^{\frac{n+1}{2}} \det(V^*/\sigma^2)^{-\frac{1}{2}} \\ &= (2\pi)^{\frac{n+1}{2}} \det(K^\top K + D\sigma^2)^{-\frac{1}{2}} \end{aligned}$$

$$\begin{aligned} \exp \left( -\frac{1}{2\sigma^2} S^* \right) &= \exp \left( -\frac{1}{2\sigma^2} \mu^\top (D\sigma^2) \mu - (\mu^*)^\top (V^*)^{-1} (\mu^*) + y^\top y \right) \\ &= \exp \left( -\frac{1}{2\sigma^2} \left( y^\top y - y^\top K (K^\top K + D\sigma^2)^{-1} K^\top y \right) \right) \\ &= \exp \left( -\frac{1}{2\sigma^2} \left( y^\top \left( I - K (K^\top K + D\sigma^2)^{-1} K^\top \right) y \right) \right) \\ &= \exp \left( -\frac{1}{2\sigma^2} \left( y^\top \left( K^\top D^{-1} K + I\sigma^2 \right)^{-1} y \right) \right) \end{aligned}$$

So

$$\begin{aligned} f(y|\lambda, \sigma^2) &= (2\pi)^{-\frac{n}{2}} \left( \frac{1}{\sigma^2} \right)^{\frac{1}{2}} \det(D)^{\frac{1}{2}} \det(K^\top K + \sigma^2 D)^{-\frac{1}{2}} \\ &\quad \times \exp \left( -\frac{1}{2\sigma^2} \left( y^\top \left( I\sigma^2 + K^\top D^{-1} K \right)^{-1} y \right) \right) \end{aligned}$$

(1) I have

$$\begin{aligned}
f(y|\sigma^2) &= \int f(y|\lambda, \sigma^2) \pi(\lambda) d\lambda \\
&= (2\pi)^{-\frac{n}{2}} (\sigma^2)^{\frac{1}{2}} \int \left[ \det(D)^{\frac{1}{2}} \right] \left[ \det(K^\top K + D\sigma^2)^{-\frac{1}{2}} \right] \\
&\quad \times \exp\left(-\frac{1}{2} \left( y^\top (I\sigma^2 + K^\top D^{-1}K)^{-1} y \right)\right) \left[ \prod_{i=0}^n \lambda_i^{a-1} \right] d\lambda_0 \dots d\lambda_n
\end{aligned}$$

because  $b = 0$ .

- It is  $\exp\left(-\frac{y^\top y}{2\sigma^2}\right) \leq \exp\left(-\frac{y^\top (I\sigma^2 + K^\top D^{-1}K)^{-1} y}{2}\right)$
- It is  $\det(D)^{\frac{1}{2}} = \prod_{i=0}^n \lambda_i^{\frac{1}{2}}$ .
- If  $\{e_j\}_{j=0}^{n-1}$  are eigenvalues of  $K^\top K$  and  $e_{\max} = \max(\{e_j\})$ , then  $K^\top K + D\sigma^2 \leq Ie_{\max} + D\sigma^2$ , consequently  $\det(K^\top K + D\sigma^2)^{-\frac{1}{2}} \geq \prod_{j=0}^n (\lambda_j \sigma^2 + e_{\max})^{-\frac{1}{2}}$ .

Then

$$\begin{aligned}
f(y|\sigma^2) &\geq (2\pi)^{-\frac{n}{2}} (\sigma^2)^{\frac{1}{2}} \int \prod_{j=0}^n \lambda_j^{\frac{1}{2}} \prod_{j=0}^n (\lambda_j \sigma^2 + e_{\max})^{-\frac{1}{2}} \exp\left(-\frac{1}{2\sigma^2} y^\top y\right) \prod_{j=0}^n \lambda_j^{a-1} d\lambda_0 \dots d\lambda_n \\
&= (2\pi)^{-\frac{n}{2}} (\sigma^2)^{\frac{1}{2}} \exp\left(-\frac{1}{2\sigma^2} y^\top y\right) \int \dots \int \prod_{j=0}^n \left[ \lambda_j^{\frac{1}{2}} \right] \left[ \prod_{j=0}^n (\lambda_j \sigma^2 + e_{\max})^{-\frac{1}{2}} \right] \left[ \prod_{j=0}^n \lambda_j^{a-1} \right] d\lambda_0 \dots d\lambda_n \\
&= (2\pi)^{-\frac{n}{2}} (\sigma^2)^{\frac{1}{2}} \exp\left(-\frac{1}{2\sigma^2} y^\top y\right) \prod_{j=0}^n \int \frac{\lambda_j^{a-\frac{1}{2}}}{(\lambda_j \sigma^2 + e_{\max})^{\frac{1}{2}}} d\lambda_j
\end{aligned}$$

Let  $t_i = 1/\lambda_i$ , then

$$f(y|\sigma^2) \geq (2\pi)^{-\frac{n}{2}} (\sigma^2)^{\frac{1}{2}} \exp\left(-\frac{1}{2\sigma^2} y^\top y\right) \prod_{j=0}^n \int \frac{t_{jj}^{-a-1}}{\left(t_j + \frac{\sigma^2}{e_{\max}}\right)^{\frac{1}{2}}} d\lambda_j$$

which is finite if and only if  $a \in (-1/2, 0)$ .

(2)

- (a) If  $a > 0$ ,  $b > 0$  then  $\lambda_i \stackrel{\text{iid}}{\sim} \text{Ga}(a, b)$  for all  $i = 1, \dots, n$ , and if  $c > 0$ ,  $d > 0$  then  $\tau \stackrel{\text{iid}}{\sim} \text{Ga}(c, d)$  which are proper. So  $\Pi(\beta, \sigma^2, \lambda, \tau)$  is a proper prior, and hence it leads to proper posterior.

(b) I have

$$\begin{aligned}
f(y|\sigma^2) &= (2\pi)^{-\frac{n}{2}} (\sigma^2)^{\frac{1}{2}} \int \left[ \det(D)^{\frac{1}{2}} \right] \left[ \det(K^\top K + D\sigma^2)^{-\frac{1}{2}} \right] \\
&\quad \times \exp\left(-\frac{1}{2} \left( y^\top (I\sigma^2 + K^\top D^{-1}K)^{-1} y \right)\right) \pi(\lambda) d\lambda
\end{aligned}$$

It is  $\det(D)^{\frac{1}{2}} = \prod_{i=0}^n \lambda_i^{\frac{1}{2}}$ . Also, it is  $K^\top K + D\sigma^2 \geq D\sigma^2$  then  $\det(K^\top K + D\sigma^2)^{-\frac{1}{2}} \leq \prod_{j=0}^n (\lambda_j \sigma^2)^{-\frac{1}{2}}$ . Hence

$$f(y|\sigma^2) \leq (2\pi)^{-\frac{n}{2}} (\sigma^2)^{\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2} y^\top (I - P) y\right) \int \pi(\lambda) d\lambda$$

which implies that  $f(y|\sigma^2) < \infty$  if  $\pi(\lambda)$  is proper. Yet,

$$\begin{aligned} f(y) &= \int f(y|\sigma^2) \pi(\sigma^2) d\sigma^2 \\ &\leq (2\pi)^{-\frac{n}{2}} \int (\sigma^2)^{-\frac{n}{2}+c+1} \exp\left(-\frac{1}{\sigma^2} \left(\frac{y^\top (I - P) y}{2} + d\right)\right) d\sigma^2 \end{aligned}$$

which is finite if  $y^\top (I - P) y + 2d > 0$  and  $c > -\frac{n}{2}$ .

(c) No. This implies  $\pi(\lambda, \sigma^2) \propto \sigma^2 \prod_{j=0}^n \lambda_j^{-1}$ . It is improper prior as  $\int \pi(\lambda, \sigma^2) d(\lambda, \sigma^2) = \infty$ , and  $(a, b, c, d) = (0, 0, 0, 0)$  which violates the necessary conditions.

(d) No, it violates the necessary conditions.

**Exercise 18.** (★) Students are encouraged to practice on the Exercises 6.1-6.19 from the textbook “Bishop, C. M. (2006). *Pattern recognition and machine learning* (Vol. 4, No. 4, p. 738). New York: Springer.” available from

- <https://www.microsoft.com/en-us/research/uploads/prod/2006/01/Bishop-Pattern-Recognition-and-Machine-Learning-2006.pdf>

The solutions are available from

- [https://blackboard.durham.ac.uk/ultra/courses/\\_44662\\_1/outline/create/document?id=\\_1396738\\_1](https://blackboard.durham.ac.uk/ultra/courses/_44662_1/outline/create/document?id=_1396738_1)



## Part 4. The kernel trick

## Part 5. Multi-class classification

## Part 6. Artificial Neural Networks

## Part 7. Gaussian process regression

## Part 8. Revision