## **Optimization in Machine Learning**

HWS 2024

Universität Mannheim Prof. Simon Weißmann, Felix Benning

## **Solution Sheet 6**

For the exercise class on the 05.12.2024.

Hand in your solutions by 10:15 in the lecture on Tuesday 03.12.2024.

## Exercise 1 (Conditional expectation).

(8 Points)

Let  $(\Omega, \mathcal{F}, \mathbb{P})$  be an underlying probability space.

- (i) Let  $\mathcal{G} \subset \mathcal{F}$  be a  $\sigma$ -algebra and let  $X, Y \in L^1(\Omega, \mathcal{F}, \mathbb{P})$ . Prove that
  - (a) for any  $\lambda \in \mathbb{R}$  there holds  $\mathbb{E}[\lambda X + Y | \mathcal{G}] = \lambda \mathbb{E}[X | \mathcal{G}] + \mathbb{E}[Y | \mathcal{G}].$

Solution. Let  $A \in \mathcal{G}$ . Then

$$\mathbb{E}\left[\mathbb{1}_{A}\left(\lambda\mathbb{E}[X \mid \mathcal{G}] + \mathbb{E}[Y \mid \mathcal{G}]\right)\right] = \lambda\mathbb{E}\left[\mathbb{1}_{A}\mathbb{E}[X \mid \mathcal{G}]\right] + \mathbb{E}\left[\mathbb{1}_{A}\mathbb{E}[Y \mid \mathcal{G}]\right]$$
$$= \lambda\mathbb{E}\left[\mathbb{1}_{A}X\right] + \mathbb{E}\left[\mathbb{1}_{A}Y\right]$$
$$= \mathbb{E}\left[\mathbb{1}_{A}(\lambda X + Y)\right]$$

Thus, by uniqueness of the conditional expectation,  $\mathbb{E}[\lambda X + Y | \mathcal{G}] = \lambda \mathbb{E}[X | \mathcal{G}] + \mathbb{E}[Y | \mathcal{G}]$ almost surely.

(b) if  $X \ge Y$   $\mathbb{P}$ -almost surely, then  $\mathbb{E}[X \mid \mathcal{G}] \ge \mathbb{E}[Y \mid \mathcal{G}]$   $\mathbb{P}$ -a.s..

Solution. Let  $A := \{\mathbb{E}[X | \mathcal{G}] < \mathbb{E}[Y | \mathcal{G}]\}$ . Then A is  $\mathcal{G}$  measurable and it follows from the definition of the conditional expectation that

$$0 \stackrel{\text{def. }A}{\geq} \mathbb{E}[\mathbb{1}_A(\mathbb{E}[X \mid \mathcal{G}] - \mathbb{E}[Y \mid \mathcal{G}])] = \mathbb{E}[\mathbb{1}_A(X - Y)] \stackrel{X \geq Y}{\geq} 0.$$

Thus,  $\mathbb{1}_A(X - Y) = 0$  almost surely and by definition of A this implies  $\mathbb{1}_A = 0$  almost surely, i.e.  $\mathbb{E}[X | \mathcal{G}] \ge \mathbb{E}[Y | \mathcal{G}]$  almost surely.  $\Box$ 

(c)  $|\mathbb{E}[X | \mathcal{G}]| \leq \mathbb{E}[|X| | \mathcal{G}].$ 

Solution. One possibility is to apply Jensen's inequality for the conditional expectation to the convex function  $x \mapsto |x|$ . Another way is to split the random variable X into its positive and negative part  $X = X^+ - X^-$  and to use the linearity of the conditional expectation and the previous result, i.e.  $X^+ \ge 0$  almost surely implies  $\mathbb{E}[X^+ | \mathcal{G}] \ge 0$  almost surely.

(3 pts)

(ii) Let  $Y_1, Y_2$  be iid. random variables with

$$\mathbb{P}(Y_1 = 2) = \mathbb{P}(Y_1 = 0.5) = \frac{1}{2},$$

and set

$$S_0 = 2, \quad S_k = S_0 \cdot \prod_{i=1}^k Y_i, \ k = 1, 2.$$

Compute  $\mathbb{E}[S_2 | \mathcal{F}_1]$ , where  $\mathcal{F}_1 = \sigma(S_1)$ .

(2 pts)

Solution. We have

$$\mathbb{E}[S_2 \mid \mathcal{F}_1] = \mathbb{E}[S_1 Y_2 \mid \mathcal{F}_1] = S_1 \mathbb{E}[Y_2 \mid \mathcal{F}_1] = S_1 \mathbb{E}[Y_2] = S_1 \left(\frac{2}{2} + \frac{1/2}{2}\right) = \frac{5}{2}S_1.$$

- (iii) Let X, Y be independent Bernoulli distributed random variables with parameter  $p \in [0, 1]$  and define  $Z := \mathbb{1}_{\{X+Y=0\}}$ .
  - (a) Compute  $\mathbb{E}[X \mid \sigma(Z)]$  and  $\mathbb{E}[Y \mid \sigma(Z)]$ .

Solution. We have

$$\mathbb{P}(Z=1) = \mathbb{P}(X+Y=0) = \mathbb{P}(X=0, Y=0) = (1-p)^2.$$
(1)

Thus,

$$\begin{split} \mathbb{E}[X \mid \sigma(Z)] &= \mathbb{E}[X \mid Z = 0] \mathbb{1}_{Z=0} + \underbrace{\mathbb{E}[X \mid Z = 1]}_{=0} \mathbb{1}_{Z=1} \\ &= \frac{\mathbb{E}[X \mathbb{1}_{Z=0}]}{\mathbb{P}(Z = 0)} \mathbb{1}_{Z=0} \\ &= \frac{0\mathbb{P}(X = 0, Y = 1) + 1[\mathbb{P}(X = 1, Y = 0) + \mathbb{P}(X = 1, Y = 1)]}{1 - (1 - p)^2} \mathbb{1}_{Z=0} \\ &= \frac{p(1 - p) + p^2}{1 - (1 - 2p + p^2)} \mathbb{1}_{Z=0} \\ &= \frac{p}{2p - p^2} \mathbb{1}_{Z=0} \\ &= \frac{1}{2 - p} \mathbb{1}_{Z=0}. \end{split}$$

By a symmetric argument we have  $\mathbb{E}[Y \mid \sigma(Z)] = \frac{1}{2-p} \mathbb{1}_{Z=0}$ .

(b) When are these random variables independent? **Hint:** You may use the fact that a real valued random variable is independent from itself if and only if it is a constant.

Solution. Since we have  $\mathbb{E}[X \mid \sigma(Z)] = \mathbb{E}[Y \mid \sigma(Z)]$  almost surely by the previous result,  $\mathbb{E}[X \mid \sigma(Z)]$  is independent from  $\mathbb{E}[Y \mid \sigma(Z)]$  if and only if it is a constant almost surely. For this to be the case, we need  $\mathbb{1}_{Z=0}$  to be constant almost surely, i.e.  $\mathbb{P}(Z=0) = 0$  or  $\mathbb{P}(Z=0) = 1$ . This is only the case for  $p \in \{0,1\}$  by (1). In summary,  $\mathbb{E}[X \mid \sigma(Z)]$  is independent from  $\mathbb{E}[Y \mid \sigma(Z)]$  if and only if  $p \in \{0,1\}$ .  $\Box$ 

(3 pts)

## Exercise 2 (Martingales).

Let  $(\Omega, \mathcal{F}, \mathbb{P})$  be an underlying probability space.

(i) Let  $Y_1, \ldots, Y_N$  be iid. random variables with  $\mathbb{E}[Y_1] = 1$  and  $\mathbb{E}[|Y_1|] < \infty$ , let  $\mathcal{F}_k := \sigma(Y_1, \ldots, Y_k)$  and define

$$S_0 = Y_0 := 1, \quad S_k = Y_1 \cdot \dots \cdot Y_k = \prod_{i=1}^{\kappa} Y_i, \ k = 1, \dots, N.$$

Prove that  $(S_k)_{k=1,\dots,N}$  is a Martingale with respect to the filtration  $(\mathcal{F}_k)_{k=1,\dots,N}$ . (2 pts)

(4 Points)

Solution. We first check  $\mathbb{E}[|S_k|] < \infty$  for all k,

$$\mathbb{E}[|S_k|] = \mathbb{E}\left[\left|\prod_{i=1}^k Y_i\right|\right] = \mathbb{E}\left[\prod_{i=1}^k |Y_i|\right] \stackrel{\text{iid}}{=} \prod_{i=1}^k \mathbb{E}[|Y_i|] < \infty$$

For the martingale property note that  $S_k$  is  $\mathcal{F}_k$ -measurable by definition, thus we have

$$\mathbb{E}[S_{k+1}|\mathcal{F}_k] = \mathbb{E}\Big[\prod_{i=1}^{k+1} Y_i \,|\, \mathcal{F}_k\Big] = \mathbb{E}\Big[Y_{k+1}S_k \,|\, \mathcal{F}_k\Big] = \mathbb{E}\Big[Y_{k+1} \,|\, \mathcal{F}_k\Big]S_k = \mathbb{E}[Y_{k+1}]S_k$$
$$= S_k.$$

(ii) Let  $(X_k)_{k\in\mathbb{N}}$  and  $(Y_k)_{k\in\mathbb{N}}$  be two martingales. Prove that  $(aX_k + bX_k)_{k\in\mathbb{N}}$  is a martingale for any  $a, b \in \mathbb{R}$ . (1 pt)

Solution. We have

$$\mathbb{E}[|aX_k + bY_k|] < |a|\mathbb{E}[|X_k|] + |b|\mathbb{E}[|Y_k|] < \infty$$

for all k and the martingale property follows from linearity of the conditional expectation

$$\mathbb{E}[aX_{k+1} + bY_{k+1}|\mathcal{F}_k] = a\mathbb{E}[X_{k+1}|\mathcal{F}_k] + b\mathbb{E}[Y_{k+1}|\mathcal{F}_k] = aX_k + bY_k.$$

(iii) Let  $(X_k)_{k\in\mathbb{N}}$  and  $(Y_k)_{k\in\mathbb{N}}$  be two super-martingales. Prove that  $(\min\{X_k, Y_k\})_{k\in\mathbb{N}}$  is a supermartingale. (1 pt)

Solution. Since  $X_k, Y_k$  are super-martingales, we have  $\mathbb{E}[|\min\{X_k, Y_k\}|] \leq \mathbb{E}[|X_k|] < \infty$ . Since min is a measurable function,  $Z_k := \min\{X_k, Y_k\}$  is  $\mathcal{F}_k$ -measurable and we obtain the super-martingale property from

$$\mathbb{E}[Z_{k+1} \mid \mathcal{F}_k] = \mathbb{E}[\min\{X_{k+1}, Y_{k+1}\} \mid \mathcal{F}_k] \le \mathbb{E}[X_{k+1} \mid \mathcal{F}_k] \le X_k$$

and similarly  $\mathbb{E}[Z_{k+1} | \mathcal{F}_k] \leq Y_k$  which together imply

$$\mathbb{E}[Z_{k+1} \mid \mathcal{F}_k] \le \min\{X_k, Y_k\}.$$

Exercise 3 (SGD with random batches).

Consider the expected risk minimization problem

$$\min_{x \in \mathbb{R}^d} F(x), \quad F(x) := \mathbb{E}_{Z \sim \mu_Z}[f(x, Z)],$$

where  $f : \mathbb{R}^d \times \mathbb{R}^p \to \mathbb{R}$  is measurable and  $Z : \Omega \to \mathbb{R}^p$  is a random vector. We assume that the "usual" conditions (conditions for Lemma 4.1.2 such as measurability, integrability, etc.) from the lecture are satisfied.

(12 Points)

(i) Assume that  $\mathbb{E}[\|\nabla_x f(x,Z) - \nabla_x F(x)\|^2] \leq c$  for some c > 0 and all  $x \in \mathbb{R}^d$ , and let  $Z^{(1)}, \ldots, Z^{(B)}, B \geq 2$ , be iid. random variables with  $Z^{(1)} \sim \mu_Z$ . Prove that

$$\mathbb{E}\left[\left\|\frac{1}{B}\sum_{m=1}^{B}\nabla_{x}f(x,Z^{(m)})-\nabla_{x}F(x)\right\|^{2}\right] \leq \frac{c}{B}$$

for all  $x \in \mathbb{R}^d$ .

Solution. From  $F(x) = \mathbb{E}[f(x, Z)]$  follows  $\nabla_x F(x) = \mathbb{E}[\nabla_x f(x, Z)]$ . Using

$$\left\|\sum_{k} v_{k}\right\|^{2} = \left\langle\sum_{k} v_{k}, \sum_{n} v_{n}\right\rangle = \sum_{k,n} \left\langle v_{k}, v_{n}\right\rangle \tag{2}$$

we have, since the  $Z^{(i)}$  are iid

$$\begin{split} & \mathbb{E}\left[\left\|\frac{1}{B}\sum_{m=1}^{B}\nabla_{x}f(x,Z^{(m)})-\nabla_{x}F(x)\right\|^{2}\right] \\ \stackrel{(2)}{=} \mathbb{E}\left[\frac{1}{B^{2}}\sum_{m,n=1}^{B}\left\langle\nabla_{x}f(x,Z^{(m)})-\nabla_{x}F(x),\nabla_{x}f(x,Z^{(n)})-\nabla_{x}F(x)\right\rangle^{2}\right] \\ &=\frac{1}{B^{2}}\sum_{m,n=1}^{B}\mathbb{E}\left[\left\langle\nabla_{x}f(x,Z^{(m)})-\nabla_{x}F(x),\nabla_{x}f(x,Z^{(n)})-\nabla_{x}F(x)\right\rangle^{2}\right] \\ \stackrel{\text{iid}}{=} \frac{1}{B^{2}}\sum_{m=1}^{B}\underbrace{\mathbb{E}\left[\left\|\nabla_{x}f(x,Z^{(m)})-\nabla_{x}F(x)\right\|^{2}\right]}_{\leq c} \\ &\leq \frac{c}{B}. \end{split}$$

Let  $X_0: \Omega \to \mathbb{R}^d$  be the initial random variable,  $(\alpha_k)_{k \in \mathbb{N}}$  be a sequence of positive deterministic step sizes and for batch sizes  $(B_k)_{k \in \mathbb{N}}$  let  $(Z_k^{(m)})_{k \in \mathbb{N}, m=1,...,B_{k-1}}$  be a sequence of iid. random variables with  $Z_1^{(1)} \sim \mu_Z$ .

(ii) Formulate the stochastic gradient descent (SGD) scheme using the stochastic gradient estimator with batch size  $B_k$ :

$$G_k(x) := \frac{1}{B_k} \sum_{m=1}^{B_k} \nabla_x f(x, Z_{k+1}^{(m)}), \quad x \in \mathbb{R}^d.$$
(2 pts)

Solution. Until convergence select

$$X_{k+1} := X_k - \alpha_k G_k(X_k)$$

and increase k (cf. Algorithm 8 in the lecture notes).

(4 pts)

(iii) Assume that F is L-smooth and  $\mu$ -strongly convex, and let  $\alpha_k \in (0, \frac{1}{L}]$  for all  $k \in \mathbb{N}$ . Prove that

$$\mathbb{E}[\|X_{k+1} - x_*\|^2] \le (1 - \alpha_k \mu) \mathbb{E}[\|X_k - x_*\|^2] + \tilde{c} \frac{\alpha_k^2}{B_k}$$

for some  $\tilde{c} > 0$ , where  $(X_k)_{k \in \mathbb{N}}$  is generated by SGD with batch-sizes  $(B_k)_{k \in \mathbb{N}}$  and  $x_* = \arg \min_{x \in \mathbb{R}^d} F(x)$ . (3 pts)

Solution. The proof is analogous to the proof of Theorem 4.1.16, i.e. let the filtration be given by  $\mathcal{F}_n := \sigma(X_0, (Z_k^{(m)})_{m=1,...,B_k}, k \leq n)$  such that by induction  $X_n$  is  $\mathcal{F}_n$ -measurable with induction step  $n \to n + 1$ 

$$X_{n+1} = X_n - \alpha_n G_n(X_n) = \underbrace{X_n}_{\mathcal{F}_n\text{-meas.}} - \alpha_n \frac{1}{B_n} \sum_{m=1}^{B_m} \underbrace{\nabla f(X_n, Z_{n+1}^{(m)})}_{\mathcal{F}_{n+1}\text{-meas.}}.$$

Let  $M_{k+1} := \nabla F(X_k) - G_k(X_k)$ , then we have

$$\begin{aligned} \|X_{k+1} - x_*\|^2 \\ &= \|X_k - \alpha_k G_k(X_k) - x_*\|^2 \\ &= \|X_k - x_*\|^2 - 2\alpha_k \langle \underbrace{G_k(X_k)}_{=\nabla F(X_k) - M_{k+1}} \rangle + \alpha_k^2 \|G_k(X_k)\|^2 \\ &= \|X_k - x_*\|^2 - 2\alpha_k \langle \nabla F(X_k), X_k - x_* \rangle + \alpha_k^2 \|\nabla F(X_k)\|^2 \\ &+ 2\alpha_k \langle M_{k+1}, X_k - x_* \rangle - 2\alpha_k^2 \langle M_{k+1}, \nabla F(X_k) \rangle + \alpha_k^2 \|M_{k+1}\|^2 \end{aligned}$$

Since  $X_k$  is  $\mathcal{F}_k$ -measurable and  $\mathbb{E}[G_k(X_k) | \mathcal{F}_k] = \nabla F(X_k)$ , we have

$$\mathbb{E}[M_{k+1} \,|\, \mathcal{F}_k] = 0 \tag{3}$$

$$\mathbb{E}\left[\left\|M_{k+1}\right\|^{2} \mid \mathcal{F}_{k}\right] \stackrel{(i)}{\leq} \frac{c}{B_{k}}.$$
(4)

This implies

$$\mathbb{E}[\|X_{k+1} - x_*\|^2 | \mathcal{F}_k] 
\stackrel{(3)}{=} \|X_k - x_*\|^2 - 2\alpha_k \langle \nabla F(X_k), X_k - x_* \rangle + \alpha_k^2 \|\nabla F(X_k)\|^2 + \alpha_k^2 \mathbb{E}[\|M_{k+1}\|^2 | \mathcal{F}_k] 
\stackrel{(4)}{\leq} \|X_k - x_*\|^2 - 2\alpha_k \langle \nabla F(X_k), X_k - x_* \rangle + \alpha_k^2 \|\nabla F(X_k)\|^2 + \alpha_k^2 \frac{c}{B_k}.$$
(5)

Recall that  $\mu$ -strong convexity implies

$$\langle \nabla F(X_k), x_* - X_k \rangle + \frac{\mu}{2} ||x_* - X_k||^2 \le F(x_*) - F(X_k)$$

which, rearranged, gives

$$-\langle \nabla F(X_k), X_k - x_* \rangle \le F(x_*) - F(X_k) - \frac{\mu}{2} \|x_* - X_k\|^2.$$
(6)

L-smoothness on the other hand implies via the descent lemma

$$F(x_*) \le F(X_k) - \frac{1}{2L} \|\nabla F(X_k)\|^2$$

which, rearranged, gives

$$F(x_*) - F(X_k) \le -\frac{1}{2L} \|\nabla F(X_k)\|^2$$
 (7)

Starting with (5) and applying (6) and then (7) we get

$$\begin{split} & \mathbb{E}\left[\left\|X_{k+1} - x_{*}\right\|^{2} | \mathcal{F}_{k}\right] \\ \stackrel{(6)}{\leq} \left\|X_{k} - x_{*}\right\|^{2} + 2\alpha_{k}\left[F(x_{*}) - F(X_{k}) - \frac{\mu}{2}\left\|x_{*} - X_{k}\right\|^{2}\right] + \alpha_{k}^{2}\left\|\nabla F(X_{k})\right\|^{2} + \alpha_{k}^{2}\frac{c}{B_{k}} \\ \stackrel{(7)}{\leq} \left\|X_{k} - x_{*}\right\|^{2} + 2\alpha_{k}\left[-\frac{1}{2L}\|\nabla F(X_{k})\|^{2} - \frac{\mu}{2}\|x_{*} - X_{k}\|^{2}\right] + \alpha_{k}^{2}\left\|\nabla F(X_{k})\right\|^{2} + \alpha_{k}^{2}\frac{c}{B_{k}} \\ &= (1 - \alpha_{k}\mu)\left\|X_{k} - x_{*}\right\|^{2} + \alpha_{k}\underbrace{(\alpha_{k} - \frac{1}{L})}_{\leq 0}\left\|\nabla F(X_{k})\right\|^{2} + \alpha_{k}^{2}\frac{c}{B_{k}} \\ &\leq (1 - \alpha_{k}\mu)\left\|X_{k} - x_{*}\right\|^{2} + c\frac{\alpha_{k}^{2}}{B_{k}}. \end{split}$$

(iv) Determine sequences of step sizes  $(\alpha_k)_{k \in \mathbb{N}}$  and batch-sizes  $(B_k)_{k \in \mathbb{N}}$  to deduce convergence  $\lim_{k \to \infty} \mathbb{E}[||X_{k+1} - x_*||^2] = 0.$  (3 pts)

Solution. There are many possible combinations of sequences, one possibility is to set the batchsize to be constant  $B_k = B$ . In this case, our upper bound is equal to the one without minibatching up to a constant and the step size from Corollary 4.1.17 works.

Another possibility is to select a constant step size  $\alpha < \frac{1}{L}$  and only change the batch size. In this case we have for  $\Delta_k = \mathbb{E}[||X_k - x_*||^2]$  by induction

$$\Delta_{n+1} \le (1 - \alpha \mu)^{n+1} \Delta_0 + c \alpha^2 \sum_{k=0}^n \frac{(1 - \alpha \mu)^{n-k}}{B_k}.$$

The induction step  $(n-1) \rightarrow n$  follows from the previous exercise, i.e.

$$\begin{split} \Delta_{n+1} &\leq (1 - \alpha \mu) \Delta_n + c \alpha^2 \frac{1}{B_n} \\ &\stackrel{\text{ind.}}{\leq} (1 - \alpha \mu) \Big[ (1 - \alpha \mu)^n \Delta_0 + c \alpha^2 \sum_{k=0}^{n-1} \frac{(1 - \alpha \mu)^{n-1-k}}{B_k} \Big] + c \alpha^2 \frac{1}{B_n} \\ &\leq (1 - \alpha \mu)^{n+1} \Delta_0 + c \alpha^2 \sum_{k=0}^n \frac{(1 - \alpha \mu)^{n-k}}{B_k}. \end{split}$$

To prove that  $0 < \alpha < \frac{1}{L}$  and  $\sum_{k=0}^{\infty} \frac{1}{B_k} < \infty$  are sufficient conditions for convergence we use

 $\frac{1}{L} \leq \frac{1}{\mu},$  which implies  $1 > (1-\alpha\mu) > 0$  and split the sum into two parts

$$\Delta_{n+1} \leq (1 - \alpha \mu)^{n+1} \Delta_0 + c \alpha^2 \Big[ \sum_{k=0}^{\lceil \frac{n}{2} \rceil} \frac{(1 - \alpha \mu)^{n - \lceil \frac{n}{2} \rceil}}{B_k} + \sum_{\substack{k=\lceil \frac{n}{2} \rceil+1}}^n \frac{(1 - \alpha \mu)^{n-k}}{B_k} \Big]$$
$$\leq \underbrace{(1 - \alpha \mu)^{n+1}}_{\to 0} \Delta_0 + c \alpha^2 \Big[ \underbrace{(1 - \alpha \mu)^{\frac{n}{2}}}_{\to 0} \sum_{\substack{k=0\\ <\infty}}^\infty \frac{1}{B_k} + \underbrace{\sum_{\substack{k=\lceil \frac{n}{2} \rceil+1}}^\infty \frac{1}{B_k}}_{\to 0} \Big]$$
$$\to 0 \qquad (n \to \infty).$$

Many other options are available and it is not clear what the optimal choice is.