






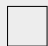








Profs. Martin Jaggi and Nicolas Flammarion
Optimization for Machine Learning – CS-439 - IC
01.07.2024 from 15h15 to 18h15
Duration : 180 minutes

Student One

SCIPER: 111111

Wait for the start of the exam before turning to the next page. This document is printed double sided, 20 pages. Do not unstaple.

- This is a closed book exam. No electronic devices of any kind.
- Place on your desk: your student ID, writing utensils, one double-sided A4 page cheat sheet if you have one; place all other personal items below your desk or on the side.
- You each have a different exam.
- For technical reasons, **do use black or blue pens for the MCQ part, no pencils!** Use white corrector if necessary.

Respectez les consignes suivantes Observe this guidelines Beachten Sie bitte die unten stehenden Richtlinien		
choisir une réponse select an answer Antwort auswählen	ne PAS choisir une réponse NOT select an answer NICHT Antwort auswählen	Corriger une réponse Correct an answer Antwort korrigieren
  		 
ce qu'il ne faut PAS faire what should NOT be done was man NICHT tun sollte		
     		



First part, multiple choice

There is **exactly one** correct answer per question.

Question 1 (Convexity & Smoothness) Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a differentiable function such that

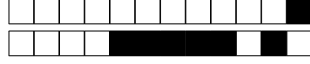
$$f(\mathbf{x}) - f(\mathbf{y}) \geq \beta \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \frac{1}{2} \gamma \|\mathbf{x} - \mathbf{y}\|^2 \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d,$$

where $\beta > 0, \gamma \geq 0$ are constants. What are all the true statements about f ?

- A) f is convex.
- B) f is Lipschitz.
- C) If $\gamma > 0$, f is strongly convex with parameter γ .
- D) None of these options.

- ☐ D
- ☐ A
- ☐ C
- ☐ B
- ☐ A and B
- ☐ A, B, and C
- ☐ B and C
- ☒ A and C

Solution: f is convex as it has a convex domain and $\beta \nabla f(\mathbf{y}) \in \partial f(\mathbf{y})$ is a subgradient, i.e., $\partial f(\mathbf{y}) \neq \emptyset$ for all $\mathbf{y} \in \mathbb{R}$. The function $g(\mathbf{x}) = f(\mathbf{x}) - \frac{1}{2} \gamma \|\mathbf{x}\|^2$ is also convex as it verifies the same condition. Therefore, f is strongly convex with parameter γ if $\gamma > 0$.



Question 2 (SGD) Consider the following optimization objective with $x \in \mathbb{R}$:

$$\min_x F(x) := \mathbb{E}_{\xi \sim \mathcal{N}(0,1)} [(x - \xi)^2].$$

A stochastic gradient descent update step is defined as

$$x_{t+1} = x_t - 2\gamma_{t+1}(x_t - \xi_t), \quad \text{where } \xi_t \sim \mathcal{N}(0, 1),$$

for all $t \geq 0$. If we run stochastic gradient descent with $x_0 = 5$ and decreasing learning rate $\gamma_t = \frac{1}{t+1}$, which of the following are true?

- A) $\lim_{t \rightarrow \infty} \mathbb{E} F(x_t) = 1$
- B) $\lim_{t \rightarrow \infty} \mathbb{E} F(x_t) = 1 + c$, where $c > 0$ is some constant dependent on x_0 .
- C) $\lim_{t \rightarrow \infty} \mathbb{E} |x_t|^2 = 0$
- D) We do not have convergence guarantee on $\mathbb{E} |x_t|^2$

What are all the true statements?

- ☐ B and D
- ☐ B and C
- ☐ A and D
- ☒ A and C

Solution: Essentially, $F(x) = x^2 - 2\mathbb{E}\xi x + \mathbb{E}\xi^2 = x^2 + 1$, which is smooth and strongly convex. The optimal solution is $x^* = 0$. SGD with diminishing step size converges to the optimal solution. (Check ex05 jupyter notebook for a visual proof)

For a formal proof, first note that C) implies A) by smoothness, i.e.,

$$F(x_t) - F(x^*) \leq \frac{L}{2} |x_t - x^*|^2.$$

Taking expectations and plugging-in $L = 2$ and $x^* = 0$ for our function,

$$\mathbb{E}[F(x_t)] - 1 \leq \mathbb{E}[|x_t|^2].$$

Taking the limit $t \rightarrow \infty$, we have $\lim_{t \rightarrow \infty} \mathbb{E} F(x_t) \leq 1$. And, since F is lower bounded by 1, we have (A).

For proving C), observe that the gradient descent step leads to the following recursion:

$$\mathbb{E}[|x_{t+1}|^2] = \mathbb{E}[|x_t - 2\gamma_{t+1}(x_t - \xi_t)|^2] = (1 - 2\gamma_{t+1})^2 \mathbb{E}[|x_t|^2] + 4\gamma_{t+1}^2.$$

Multiplying with $1/\gamma_{t+1}$ and summing over t , we have

$$\sum_{t=0}^{T-1} \frac{1}{\gamma_{t+1}} \mathbb{E}[|x_{t+1}|^2] = \sum_{t=0}^{T-1} \frac{(1 - 2\gamma_{t+1})^2}{\gamma_{t+1}} \mathbb{E}[|x_t|^2] + \sum_{t=0}^{T-1} 4\gamma_{t+1}.$$

Putting together the telescoping sum, we have

$$\frac{(1 - 2\gamma_{T+1})^2}{\gamma_{T+1}} \mathbb{E}[|x_T|^2] + \sum_{t=1}^T \left(\frac{1}{\gamma_t} - \frac{(1 - 2\gamma_t)^2}{\gamma_{t+1}} \right) \mathbb{E}[|x_t|^2] = \sum_{t=1}^T 4\gamma_t.$$

Our selection of γ_t satisfies $\frac{1}{\gamma_t} - \frac{(1 - 2\gamma_{t+1})^2}{\gamma_{t+1}} \geq 0$ for all $t \geq 0$, and hence

$$\mathbb{E}[|x_T|^2] \leq \frac{4\gamma_{T+1}}{(1 - 2\gamma_{T+1})^2} \sum_{t=1}^T \gamma_t.$$

Taking the limit $t \rightarrow \infty$, we have 0 on the RHS.

For your examination, preferably print documents compiled from auto-multiple-choice.



Question 3 (Convergence Rates) Consider the following function:

$$f(x) = \begin{cases} |x|^3 & \text{for } |x| \leq 1, \\ 3|x| - 2 & \text{otherwise.} \end{cases}$$

with domain $D_f = \mathbb{R}$. Assume we want to find a point x_T with $f(x_T) \leq \varepsilon$ starting from an unknown $x_0 \in \mathbb{R}$. Which algorithm provides the tightest applicable bound (assuming appropriate hyperparameters)?

- ☐ Gradient Descent, $T \in \mathcal{O}(\log(\frac{1}{\varepsilon}))$
- ☒ Nesterov acceleration, $T \in \mathcal{O}(\frac{1}{\sqrt{\varepsilon}})$
- ☐ Gradient Descent, $T \in \mathcal{O}(\frac{1}{\varepsilon^2})$
- ☐ Newton's method, $T \in \mathcal{O}(\log \log \frac{1}{\varepsilon})$
- ☐ Gradient Descent, $T \in \mathcal{O}(\frac{1}{\varepsilon})$

Solution: The function does not have an invertible Hessian so the standard Newton's method is not applicable. The function is not strongly convex due to the linear portions, so the stronger gradient descent bound does not hold. The function is smooth and convex, so the Nesterov acceleration bound is the fastest applicable one.

Question 4 (PowerSGD) PowerSGD approximates a matrix $M \in \mathbb{R}^{m \times n}$ with PQ^\top where $P \in \mathbb{R}^{n \times r}$ and $Q \in \mathbb{R}^{m \times r}$. What is the amount of data transmitted in PowerSGD relative to (divided by) that of standard distributed SGD?

- ☐ 1
- ☒ $\frac{mr+nr}{mn}$
- ☐ $\frac{m+n+r^2}{mn}$
- ☐ $\frac{r}{m+n}$
- ☐ $\frac{r}{mn}$
- ☐ $\frac{m+n-r}{m+n}$
- ☐ $\frac{1}{r}$

Solution: PowerSGD sends and receives P and Q instead of M . Depending on the implementation, each communication round can require several partial sums of each matrix to be transmitted, but PowerSGD does not affect this factor. Overall the communication of PowerSGD is therefore proportional to the number of elements in P and Q , i.e. $nr + mr$, and for standard training this is proportional to the number of elements in M , i.e. mn . This gives a ratio of $\frac{mr+nr}{mn}$.



Question 5 (Projections) Let $X \subseteq \mathbb{R}^d$ be closed and convex, $\mathbf{x} \in X, \mathbf{y} \in \mathbb{R}^d, \mathbf{z} \in \mathbb{R}^d$ and Π_X be the projection operator onto X . Which of the following statements are true?

- A) $\|\mathbf{x} - \Pi_X(\mathbf{y})\|^2 + \|\mathbf{x} - \Pi_X(\mathbf{z})\|^2 \leq \|2\mathbf{x} - (\mathbf{y} + \mathbf{z})\|^2$
- B) $\|\mathbf{y} - \Pi_X(\mathbf{z})\| + \|\mathbf{x} - \Pi_X(\mathbf{z})\| \geq \|\mathbf{x} - \mathbf{y}\|$
- C) $(\mathbf{x} - \Pi_X(\mathbf{y}))^\top (\mathbf{y} - \Pi_X(\mathbf{y})) \leq 0$
- D) $(\mathbf{x} - \Pi_X(\mathbf{x}))^\top (\mathbf{y} - \Pi_X(\mathbf{z})) \leq 0$
- E) $\|\mathbf{x} - \Pi_X(\mathbf{y})\| + \|\mathbf{x} - \Pi_X(\mathbf{z})\| \leq \|2\mathbf{x} - (\mathbf{y} + \mathbf{z})\|$

- ☐ A, D, and E
- ☒ B, C, and D
- ☐ B, C, and E
- ☐ A, B, and D
- ☐ A, C, and E

Solution: A counterexample for (A): let $x = 0$ be the center of a circle which is X and choose y and z outside of X such that $x = y + z$. Then, the right-hand side is zero and the left-hand side is positive. Proof for (B): Triangle inequality. Proof for (C): A fundamental fact about the projection operator. Proof for (D): $\mathbf{x} - \Pi_X(\mathbf{x}) = 0$ A counterexample for (E): the same counterexample as in (A).

Question 6 (Non-convex optimization) Let f be a non-convex function on \mathbb{R}^d that is lower bounded by some constant $B \in \mathbb{R}$, i.e., $f(\mathbf{x}) \geq B$ for all $\mathbf{x} \in \mathbb{R}^d$. Further, assume that f is smooth on a set $X \subset \mathbb{R}^d$ with parameter L . Let x_1, \dots, x_T be the trajectory obtained by running gradient descent with step size $\gamma = \frac{1}{L}$ from a $x_0 \in X$. Assume that for all x_0 and T , the trajectory stays within X . Which of the following statements are necessarily true?

- A) $\lim_{T \rightarrow \infty} f(\mathbf{x}_T) = f(\mathbf{x}^*)$ for some local minimum \mathbf{x}^* .
- B) $\lim_{T \rightarrow \infty} \|\nabla f(\mathbf{x}_T)\| = 0$.
- C) f has the sufficient decrease property in X .
- D) The trajectory x_0, x_1, \dots, x_T converges to a critical point.

- ☐ A and C
- ☐ A and B
- ☐ B and D
- ☐ A and D
- ☐ C and D
- ☒ B and C

Solution: (A) there might be no local minimum (e.g., $f(x) = e^{-x}$). (B) see Theorem 6.2. in the lecture notes or slide 29 of Lecture 5. (C) see Lemma 2.7. in the lecture notes or slide 14 of Lecture 2. (D) there could be divergent behaviors (e.g., $f(x) = e^{-x}$).



Question 7 Let f be a L -smooth convex function and recall the Nesterov acceleration algorithm given by

$$\begin{aligned}\mathbf{y}_{t+1} &:= \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t), \\ \mathbf{z}_{t+1} &:= \mathbf{z}_t - \frac{t+1}{2L} \nabla f(\mathbf{x}_t), \\ \mathbf{x}_{t+1} &:= \frac{t+1}{t+3} \mathbf{y}_{t+1} + \frac{2}{t+3} \mathbf{z}_{t+1}.\end{aligned}$$

Which of the following statements are true ?

- A) $f(\mathbf{y}_{t+1}) \leq f(\mathbf{x}_t)$.
B) $\|\mathbf{y}_{t+1} - \mathbf{x}_*\|^2 \leq \|\mathbf{x}_t - \mathbf{x}_*\|^2$, where \mathbf{x}_* is a global minimum.
C) For any $t \geq 0$, we have $\|\mathbf{z}_t - \mathbf{x}_*\| \leq \|\mathbf{z}_0 - \mathbf{x}_*\|$, where \mathbf{z}_0 is the initialization.

☒ All of the above

☐ B and C

☐ A

☐ A and B

☐ B

☐ C

☐ A and C

Solution: D can be inferred from the potential function used in Theorem 2.9. B and C are direct consequence of the descent lemma. Let $\mathbf{x}_+ = \mathbf{x} - \frac{1}{L} \nabla f(\mathbf{x})$. Using the smoothness property of f , we have,

$$f(\mathbf{x}_+) \leq f(\mathbf{x}) + (\nabla f(\mathbf{x}))^\top (\mathbf{x}_+ - \mathbf{x}) + \frac{L}{2} \|\mathbf{x}_+ - \mathbf{x}\|^2 = f(\mathbf{x}) - \frac{1}{2L} \|\nabla f(\mathbf{x})\|^2.$$

Hence, we have,

$$\|\nabla f(\mathbf{x})\|^2 \leq 2L(f(\mathbf{x}) - f(\mathbf{x}_+)) \leq 2L(f(\mathbf{x}) - f(\mathbf{x}_*))$$

From the first part of inequality, we have, $f(\mathbf{x}_+) \leq f(\mathbf{x})$. Coming to the distance,

$$\|\mathbf{x}_+ - \mathbf{x}_*\|^2 = \|\mathbf{x} - \mathbf{x}_* - \frac{1}{L} \nabla f(\mathbf{x})\|^2 = \|\mathbf{x} - \mathbf{x}_*\|^2 + \frac{1}{L^2} \|\nabla f(\mathbf{x})\|^2 - \frac{2}{L} (\nabla f(\mathbf{x}))^\top (\mathbf{x} - \mathbf{x}_*).$$

Using $\|\nabla f(\mathbf{x})\|^2 \leq 2L(f(\mathbf{x}) - f(\mathbf{x}_*))$, we have

$$\|\mathbf{x}_+ - \mathbf{x}_*\|^2 \leq \|\mathbf{x} - \mathbf{x}_*\|^2 + \frac{2}{L} [f(\mathbf{x}) + (\nabla f(\mathbf{x}))^\top (\mathbf{x}_* - \mathbf{x}) - f(\mathbf{x}_*)].$$

Using the convexity of f , we have, $f(\mathbf{x}_*) \geq f(\mathbf{x}) + (\nabla f(\mathbf{x}))^\top (\mathbf{x}_* - \mathbf{x})$. Substituting this we have, $\|\mathbf{x}_+ - \mathbf{x}_*\|^2 \leq \|\mathbf{x} - \mathbf{x}_*\|^2$.



Question 8 (Subgradients) Let $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^d$ be vectors such that none of them can be written as a linear combination of the others. Let X be the convex hull of this point set, i.e.:

$$X = \left\{ \mathbf{x} \in \mathbb{R}^d \mid \mathbf{x} = \sum_{i=1}^N \alpha_i \mathbf{x}_i, \sum_{i=1}^N \alpha_i \leq 1, \alpha_i \geq 0 \right\}.$$

Let $f : X \rightarrow \mathbb{R}$ be a function such that

$$f(\mathbf{x}) = \sum_{i=1}^N \beta_i b_i \quad \text{for } \mathbf{x} = \sum_{i=1}^N \beta_i \mathbf{x}_i \in X,$$

where $b_1, \dots, b_N \in \mathbb{R}$ are fixed scalars. Which of the following statements are necessarily true?

- A) f is strongly convex.
- B) $(b_1, \dots, b_N) \in \partial f(x)$ for all x .
- C) f has exactly one subgradient at all points.
- D) None of these options.

☒ D

☐ A, B and C

☐ C

☐ B

☐ A and B

☐ A and C

☐ B and C

☐ A

Solution: (A) f is not strongly convex as it is a linear combination of non-strongly convex (linear) functions. (B) Consider $\mathbf{x}_1 = e_2, \mathbf{x}_2 = e_1$ and $b_1 = 2, b_2 = 3$. Then, $f(\mathbf{x}) = 3x_1 + 2x_2$. The only subgradient at $\mathbf{x} = (0.5, 0.5)$ is $\{(3, 2)\}$ not $\{(2, 3)\}$. (C) f can have multiple subgradients at the boundary of X .



Second part, true/false questions

Question 9 (Smoothness) Any twice differentiable function with bounded Hessians over some convex set X is smooth over X .

☒ TRUE ☐ FALSE

Solution: Yes, see Lemma 6.1 of the lecture notes.

Question 10 (Convexity) Let f be a function on a convex set $X = \cup_{i=1}^n X_i$ for some fixed $n > 0$ and X_i are also convex. If f is convex on each X_i , then f is convex on X .

☐ TRUE ☒ FALSE

Solution: Two one-dimensional square functions glued side by side is an easy counterexample.

Question 11 (Convexity) Let f be a function on a convex set $X = \cup_{i=1}^n X_i$ for some fixed $n > 0$ and X_i are also convex. If f is strictly convex on each X_i , then there are at most n local minimums of f on X .

☒ TRUE ☐ FALSE

Solution: Each X_i has at most one local minimum.

Question 12 (Stochastic Gradient Descent) If we perform stochastic gradient descent on a L -smooth function using step size $1/L$, assuming unbiased stochastic gradients, each update step guarantees a non-negative decrease in the training loss in expectation.

☐ TRUE ☒ FALSE

Solution: With an unbiased stochastic gradient, we have $\mathbb{E}g(x_t) = \nabla f(x_t)$. We have two consecutive iterates as follows:

$$\mathbb{E}_t[f(x_{t+1})] \leq f(x_t) - \eta \|\nabla f(x_t)\|_2^2 + \frac{L}{2} \eta^2 \mathbb{E}_t[\|g(x_t)\|_2^2]$$

With $\eta = 1/L$, a sufficient decrease in the expected training loss is not guaranteed, due to the variance in the stochastic gradients, as $\mathbb{E}_t[\|g(x_t)\|_2^2] \neq \|\nabla f(x_t)\|_2^2$.

Question 13 (Non-Convex Optimization) Running stochastic gradient descent to minimize non-convex functions that are bounded from below guarantees convergence to a local minimum.

☐ TRUE ☒ FALSE

Solution: Can as well be a saddle point or a region of very flat but nonzero gradients.

Question 14 (Newton) Running Newton's method on $\cos(x)$ for $x \in \mathbb{R}$ will result in convergence to a local minimum $(2k+1)\pi$ for some integer k .

☐ TRUE ☒ FALSE

Solution: Newton's method can also converge to a local maximum depending on the starting point or be undefined where the curvature is zero.



Question 15 (SignSGD) Let $f(\mathbf{x}) := \mathbf{x}^\top \mathbf{u}$ with $\mathbf{x}, \mathbf{u} \in \mathbb{R}^d$. A SignSGD step is defined by $\mathbf{x}_{t+1} = \mathbf{x}_t - \gamma \cdot \text{sign}(\mathbf{g}_t)$ where $\mathbf{g}_t = \nabla f(\mathbf{x}_t)$ and the operations are performed element-wise. Claim: There exists a $\mathbf{u} \in \mathbb{R}^d$ and $\gamma \geq 0$ such that a SignSGD update results in $f(\mathbf{x}_{t+1}) > f(\mathbf{x}_t)$.

☐ TRUE ☒ FALSE

Solution: The gradient is $\mathbf{g}_t = \mathbf{u}$ and $f(x_{t+1}) = f(x_t - \gamma \cdot \text{sign}(\mathbf{g}_t)) = \langle x_t - \gamma \cdot \text{sign}(\mathbf{g}_t), \mathbf{u} \rangle = f(x_t) - \gamma \cdot \langle \text{sign}(\mathbf{u}), \mathbf{u} \rangle$. Now, $\langle \text{sign}(\mathbf{u}), \mathbf{u} \rangle = \sum_i \text{sign}(u_i) u_i = \sum_i |u_i| \geq 0$, giving $f(x_{t+1}) \leq f(x_t)$ for all $\mathbf{u} \in \mathbb{R}^d$ and $\gamma \geq 0$.

Question 16 (Adam) An Adam update step can be formulated as:

$$\begin{aligned}\mathbf{m}_t &= \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t \\ \mathbf{v}_t &= \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \mathbf{g}_t^2 \\ \mathbf{x}_t &= \mathbf{x}_{t-1} - \gamma \cdot \frac{\mathbf{m}_t / (1 - \beta_1^t)}{\sqrt{\mathbf{v}_t / (1 - \beta_2^t)} + \varepsilon}\end{aligned}$$

where all operations are performed element-wise, $\mathbf{x}_t \in \mathbb{R}^d$ is the parameter vector, $\mathbf{m}_t \in \mathbb{R}^d$ a momentum vector, $\mathbf{v}_t \in \mathbb{R}^d$ a vector of second moment estimates, and $\mathbf{g}_t \in \mathbb{R}^d$ the gradients w.r.t. \mathbf{x}_t . The momentum and second moment vectors are initialized to $\mathbf{0}$ at $t = 0$ and we only consider updates for $t \geq 1$.

Claim: In the limit $\varepsilon \rightarrow 0$, we can set the values of β_1, β_2 in a way that recovers the SignSGD update.

☒ TRUE ☐ FALSE

Solution: Set $\beta_1 = \beta_2 = 0$ giving $\mathbf{m}_t = \mathbf{g}_t$, $\mathbf{v}_t = \mathbf{g}_t^2$, and bias correction terms of $(1 - \beta_1^t) = 1$ and $(1 - \beta_2^t) = 1$. Then $\frac{\mathbf{m}_t}{\sqrt{\mathbf{v}_t} + \varepsilon} = \frac{\mathbf{g}_t}{|\mathbf{g}_t| + \varepsilon} \rightarrow \text{sign}(\mathbf{g}_t)$ when $\varepsilon \rightarrow 0$.

Question 17 (Federated Learning) In Federated Learning each worker shares their data with the coordinating server, but not other workers.

☐ TRUE ☒ FALSE

Solution: In Federated Learning workers share model updates with the coordinating server, but their data is meant to stay private and is not shared directly.

Question 18 (Adversarial Examples) The creation of adversarial examples for neural networks can be formulated as a constrained optimization problem.

☒ TRUE ☐ FALSE

Solution: Yes, see Lecture 12.

Question 19 (Projected Gradient Descent) Projected gradient descent on the set,

$$B_1(R) := \left\{ \mathbf{x} \in \mathbb{R}^d \mid \|\mathbf{x}\|_1 = \sum_{i=1}^d |x_i| \leq R \right\},$$

can be implemented at the same time complexity (in \mathcal{O}) up to log factors as unconstrained gradient descent.

Hint: i.e., they have the same dependency on dimension d up to constants and log terms.

☒ TRUE ☐ FALSE

Solution: See Theorem 3.13. of the lecture notes and comment below or the last slide of Lecture 3.



Question 20 (Proximal Gradient Descent) Let $g, h : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex and smooth with parameter L and let $f(\mathbf{x}) = g(\mathbf{x}) + h(\mathbf{x})$. Considering minimizing over f with proximal gradient descent given proximal mapping for h and classical gradient descent for appropriately chosen stepsizes. The error as a function of the number of steps T scales the same for classical gradient descent and proximal gradient descent up to constants.

☒ TRUE ☐ FALSE

Solution: Both scales as $1/T$.

Question 21 (Proximal Gradient Descent) A step of proximal gradient descent on the function $f(\mathbf{x}) + \frac{1}{2}\|\mathbf{x}\|_2^2$ is equivalent to a step of gradient descent on the function $f(\mathbf{x})$ for properly chosen stepsizes.

☐ TRUE ☒ FALSE

Solution: different algorithms

DRAFT



Third part, open questions

Answer in the space provided! Your answer must be justified with all steps. Do not cross any checkboxes, they are reserved for correction.

For the following, assume f is a function with a L -Lipschitz continuous gradient. Let μ be some positive constant (*not necessarily the same for every question*). We focus on a basic unconstrained optimization problem

$$\operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) .$$

Let \mathbf{x}_p be the projection of \mathbf{x} onto the solution set \mathcal{X}^* , here we assume that the projection is well-defined. Let $f^* = f(\mathbf{x})$, where $\mathbf{x} \in \mathcal{X}^*$.

In the following few questions, we will study the convergence of standard iterative methods such as gradient decent and coordinate descent. We will slightly change the viewpoint and consider assumptions on the objective functions slightly different from the ones seen in class. The first few exercises in this part will be to better understand these assumptions, before then moving to the algorithm convergence analysis.

Relationship between conditions

We start by recalling the definition of strong convexity.

Strong Convexity. For any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, we have

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|^2 . \quad (\text{SC})$$

Question 22: 1 point. Show that strong convexity (SC) implies weak strong convexity (WSC).

Weak Strong Convexity. For any $\mathbf{x} \in \mathbb{R}^d$, we have

$$f^* \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}_p - \mathbf{x} \rangle + \frac{\mu}{2} \|\mathbf{x}_p - \mathbf{x}\|^2 . \quad (\text{WSC})$$



Solution: If we let $\mathbf{y} = \mathbf{x}_p$ in the SC inequality, we have the WSC inequality immediately.

Question 23: 1 point.

Show that weak strong convexity (WSC) implies the restricted secant inequality (RSI).

Restricted Secant Inequality. For any $\mathbf{x} \in \mathbb{R}^d$, we have

$$\langle \nabla f(\mathbf{x}), \mathbf{x} - \mathbf{x}_p \rangle \geq \frac{\mu}{2} \|\mathbf{x} - \mathbf{x}_p\|^2 . \quad (\text{RSI})$$



Solution: Re-arrange the WSC inequality to

$$\langle \nabla f(\mathbf{x}), \mathbf{x} - \mathbf{x}_p \rangle \geq f(\mathbf{x}) - f^* + \frac{\mu}{2} \|\mathbf{x}_p - \mathbf{x}\|^2 .$$

Since $f(\mathbf{x}) - f^* \geq 0$, we have RSI with $\frac{\mu}{2}$.

Question 24: 1 point.

Show that the restricted secant inequality (RSI) implies the following error bound (EB).

For your examination, preferably print documents compiled from auto-multiple-choice.



Error Bound. For any $\mathbf{x} \in \mathbb{R}^d$, we have,

$$\|\nabla f(\mathbf{x})\| \geq \mu \|\mathbf{x} - \mathbf{x}_p\|. \quad (\text{EB})$$

☐ 0 ☒ 1

Solution: Using Cauchy-Schwartz on the RSI we have

$$\|\nabla f(\mathbf{x})\| \|\mathbf{x} - \mathbf{x}_p\| \geq \langle \nabla f(\mathbf{x}), \mathbf{x} - \mathbf{x}_p \rangle \geq \frac{\mu}{2} \|\mathbf{x}_p - \mathbf{x}\|^2.$$

Question 25: 2 points. Assume that PL inequality (PL) implies quadratic growth (QG). Show that the above error bound (EB) is *equivalent* to having the PL inequality.

Polyak-Lojasiewicz Inequality. For any $\mathbf{x} \in \mathbb{R}^d$, we have,

$$\frac{1}{2} \|\nabla f(\mathbf{x})\|^2 \geq \mu(f(\mathbf{x}) - f^*). \quad (\text{PL})$$

Quadratic Growth. For any $\mathbf{x} \in \mathbb{R}^d$, we have,

$$f(\mathbf{x}) - f^* \geq \frac{\mu}{2} \|\mathbf{x}_p - \mathbf{x}\|^2. \quad (\text{QG})$$

☐ 0 ☐ 1 ☒ 2

Solution:

EB \rightarrow PL: By Lipschitz continuity we have

$$f(\mathbf{x}) \leq f(\mathbf{x}_p) + \langle \nabla f(\mathbf{x}_p), \mathbf{x} - \mathbf{x}_p \rangle + \frac{L}{2} \|\mathbf{x}_p - \mathbf{x}\|^2$$

use EB along with $f(\mathbf{x}_p) = f^*$ and $\nabla f(\mathbf{x}_p) = 0$, we have

$$f(\mathbf{x}) - f^* \leq \frac{L}{2} \|\mathbf{x}_p - \mathbf{x}\|^2 \leq \frac{L}{2\mu} \|\nabla f(\mathbf{x})\|^2$$

PL holds for constant $\frac{\mu}{L}$

PL \rightarrow EB: As PL implies QG, we directly have

$$\frac{1}{2} \|\nabla f(\mathbf{x})\|^2 \geq \mu(f(\mathbf{x}) - f^*) \geq \frac{\mu^2}{2} \|\mathbf{x} - \mathbf{x}_p\|^2$$

In the following, we are going to show that the PL inequality (PL) implies the quadratic growth (QG) condition. To be able to prove this statement, we break the proof into substeps.

Question 26: 1 point. Construct $g(\mathbf{x}) := \sqrt{f(\mathbf{x}) - f^*}$, show that if PL holds for f , $\|\nabla g(\mathbf{x})\|$ is lower bounded for $\mathbf{x} \notin \mathcal{X}^*$ as follows

$$\|\nabla g(\mathbf{x})\| \geq \sqrt{\frac{\mu}{2}}$$

☐ 0 ☒ 1

Solution: If we assume that f satisfies the PL inequality then for any $\mathbf{x} \notin \mathcal{X}^*$ we have

$$\|\nabla g(\mathbf{x})\|^2 = \left\| \frac{1}{2\sqrt{f(\mathbf{x}) - f^*}} \nabla f(\mathbf{x}) \right\|^2 = \frac{\|\nabla f(\mathbf{x})\|^2}{4(f(\mathbf{x}) - f^*)} \geq \frac{\mu}{2}$$

That is

$$\|\nabla g(\mathbf{x})\| \geq \sqrt{\frac{\mu}{2}}$$



Question 27: 3 points. Note that by construction, we have $g(\mathbf{x}) := \sqrt{f(\mathbf{x}) - f^*} \geq 0$. Recall the optimal solution set of f , \mathcal{X}^* such that for all $\mathbf{y} \in \mathcal{X}^*$, $g(\mathbf{y}) = 0$. Define $\mathbf{x}(t)$, for $t \geq 0$, as the solution of the differential equation, for $\mathbf{x}(t) \notin \mathcal{X}^*$,

$$\frac{d\mathbf{x}(t)}{dt} = -\nabla g(\mathbf{x}(t)),$$

$$\mathbf{x}(0) = \mathbf{x}_0.$$

Show that $\mathbf{x}(t)$ satisfies the following property,

$$g(\mathbf{x}_0) - g(\mathbf{x}(t)) \geq \frac{\mu}{2}t.$$

Furthermore, let T be the time the curve $\mathbf{x}(t)$ hits \mathcal{X}^* , i.e., $\lim_{t \rightarrow T} \mathbf{x}(t) \in \mathcal{X}^*$. What can be said about T , is it bounded ?

☐ 0 ☐ 1 ☐ 2 ☒ 3

Solution: For any point $\mathbf{x}_0 \notin \mathcal{X}^*$, consider solving the following differential equation:

$$\frac{d\mathbf{x}(t)}{dt} = -\nabla g(\mathbf{x}(t))$$

$$\mathbf{x}(t=0) = \mathbf{x}_0$$

for $\mathbf{x}(t) \notin \mathcal{X}^*$. (This is a flow orbit starting at \mathbf{x}_0 and flowing along the gradient of g .) By Question 26, ∇g is bounded from below, and by construction g is also bounded from below. Thus, by moving along the path defined by $\mathbf{x}(t)$ we are sufficiently reducing the function and will eventually reach the optimal set. Thus there exists a T such that $\mathbf{x}(T) \in \mathcal{X}^*$ (and at this point the differential equation ceases to be defined). We can show this by using the steps

$$\begin{aligned} g(\mathbf{x}_0) - g(\mathbf{x}_t) &= \int_{\mathbf{x}_t}^{\mathbf{x}_0} \langle \nabla g(\mathbf{x}), d\mathbf{x} \rangle \\ &= - \int_{\mathbf{x}_0}^{\mathbf{x}_t} \langle \nabla g(\mathbf{x}), d\mathbf{x} \rangle \\ &= - \int_0^T \langle \nabla g(\mathbf{x}(t)), \frac{d\mathbf{x}(t)}{dt} \rangle dt \\ &= \int_0^T \|\nabla g(\mathbf{x}(t))\|^2 dt \\ &\geq \int_0^T \frac{\mu}{2} dt \\ &= \frac{\mu}{2} T \end{aligned}$$

As $g(\mathbf{x}_t) \geq 0$, this shows we need to have $T \leq 2g(\mathbf{x}_0)/\mu$.

Question 28: 2 points. The length of the any curve $(\mathbf{x}(t))_{t \geq 0}$ from time t_1 to t_2 is defined as,

$$\mathcal{L}(\mathbf{x}(t), t_1, t_2) = \int_{t_1}^{t_2} \left\| \frac{d\mathbf{x}(t)}{dt} \right\| dt.$$

Let T be the time the curve hits \mathcal{X}^* . Using the definition of the length of the curve, show that

$$\int_0^T \|\nabla g(\mathbf{x}(t))\| dt \geq \|\mathbf{x}_0 - \mathbf{x}_p\|$$

☐ 0 ☐ 1 ☒ 2

Solution: The length of the orbit $\mathbf{x}(t)$ starting at \mathbf{x}_0 , which we'll denote by $\mathcal{L}(\mathbf{x}_0)$, is given by

$$\mathcal{L}(\mathbf{x}_0) = \int_0^T \left\| \frac{d\mathbf{x}(t)}{dt} \right\| dt = \int_0^T \|\nabla g(\mathbf{x}(t))\| dt \geq \|\mathbf{x}_0 - \mathbf{x}_p\|$$

the inequality follows because the orbit is a path from \mathbf{x}_0 to a point in \mathcal{X}^* (and thus it must be at least as long as the projection distance).



Question 29: 2 points. Combine the results from the previous questions (no matter if you have solved them or not), to finally prove that

$$g(\mathbf{x}_0) \geq \sqrt{\frac{\mu}{2}} \|\mathbf{x}_0 - \mathbf{x}_p\|.$$

As a consequence of this, we have that the quadratic growth condition (QG) holds directly.

☐ 0 ☐ 1 ☒ 2

Solution:

$$\begin{aligned} g(\mathbf{x}_0) - g(\mathbf{x}_T) &= \int_0^T \|\nabla g(\mathbf{x}(t))\|^2 dt \\ &\geq \sqrt{\frac{\mu}{2}} \int_0^T \|\nabla g(\mathbf{x}(t))\| dt \\ &\geq \sqrt{\frac{\mu}{2}} \|\mathbf{x}_0 - \mathbf{x}_p\| \end{aligned}$$

As $g(\mathbf{x}_T) = 0$, this yields the desired result.

Question 30: 1 point. Until now, we can summarize that we have shown that

$$SC \Rightarrow WSC \Rightarrow RSI \Rightarrow EB \equiv PL \Rightarrow QG.$$

If we further assume f is convex, we can even claim that $RSI \equiv EB \equiv PL \equiv QG$. Can you prove this statement by showing QG implies RSI?

☐ 0 ☒ 1

Solution: If we assume QG, we have

$$\frac{\mu}{2} \|\mathbf{x}_p - \mathbf{x}\|^2 \leq f(\mathbf{x}) - f^* \stackrel{\text{convexity of } f}{\leq} \langle \nabla f(\mathbf{x}), \mathbf{x} - \mathbf{x}_p \rangle$$

Convergence Proof for Gradient and Coordinate Descent via the PL Inequality

Gradient Descent

Question 31: 2 points. Consider the function f with L -Lipschitz continuous gradient, a non-empty solution set \mathcal{X}^* , and satisfies the PL inequality with constant μ . Show that gradient method with a stepsize of $1/L$

$$\mathbf{x}_{t+1} := \mathbf{x}_t - \frac{1}{L} \nabla f(\mathbf{x}_t)$$

has a global linear convergence rate, i.e.

$$f(\mathbf{x}_t) - f^* \leq \left(1 - \frac{\mu}{L}\right)^t (f(\mathbf{x}_0) - f^*).$$

☐ 0 ☐ 1 ☒ 2

Solution: If we plug the gradient update step into the Lipschitz inequality condition, we have

$$f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) \leq -\frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2$$

Now use the PL inequality, we get

$$f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) \leq -\frac{\mu}{L} (f(\mathbf{x}_t) - f^*)$$

Re-arrange the above inequality we have

$$f(\mathbf{x}_{t+1}) - f^* \leq \left(1 - \frac{\mu}{L}\right) (f(\mathbf{x}_t) - f^*)$$

Applying this inequality recursively gives the result.



Randomized Coordinate Descent

Question 32: 3 points.

Now, consider f with *coordinate-wise* L -Lipschitz-continuous gradient, i.e., for every coordinate $i = 1, \dots, d$, we have,

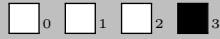
$$f(\mathbf{x} + \alpha \mathbf{e}_i) \leq f(\mathbf{x}) + \alpha \nabla_i f(\mathbf{x}) + \frac{L}{2} \alpha^2.$$

As before, we continue to assume f has a non-empty solution set \mathcal{X}^* , and satisfies the PL inequality with constant μ . Show that the coordinate descent method with a stepsize of $1/L$

$$\mathbf{x}_{t+1} := \mathbf{x}_t - \frac{1}{L} \nabla_{i_t} f(\mathbf{x}_t) \mathbf{e}_{i_t},$$

has an expected linear convergence rate if we choose the variable to update i_t uniformly at random, i.e.

$$\mathbb{E}[f(\mathbf{x}_t) - f^*] \leq \left(1 - \frac{\mu}{dL}\right)^t [f(\mathbf{x}_0) - f^*].$$



Solution:

By using the update rule in the Lipschitz condition, we have

$$f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) \leq -\frac{1}{2L} |\nabla_{i_t} f(\mathbf{x}_t)|^2$$

Taking expectation over i_t , we have

$$\begin{aligned} \mathbb{E}[f(\mathbf{x}_{t+1})] &\leq f(\mathbf{x}_t) - \frac{1}{2L} \mathbb{E}[|\nabla_{i_t} f(\mathbf{x}_t)|^2] \\ &= f(\mathbf{x}_t) - \frac{1}{2L} \sum_i \frac{1}{d} |\nabla_i f(\mathbf{x}_t)|^2 \\ &= f(\mathbf{x}_t) - \frac{1}{2dL} \|\nabla f(\mathbf{x}_t)\|^2 \end{aligned}$$

Applying PL inequality and re-arrange, we have

$$\mathbb{E}[f(\mathbf{x}_{t+1}) - f^*] \leq \left(1 - \frac{\mu}{dL}\right) \mathbb{E}[f(\mathbf{x}_t) - f^*]$$

Greedy Coordinate Descent

Question 33: 3 points.

Let f satisfy the same conditions as in the last question, but now we sample i_t according to the rule $i_t = \operatorname{argmax}_j [\nabla_j f(\mathbf{x}_t)]$. Further, assume f satisfies

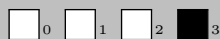
$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\mu_1}{2} \|\mathbf{y} - \mathbf{x}\|_1^2,$$

which leads to the PL inequality in the ∞ -norm:

$$\frac{1}{2} \|\nabla f(\mathbf{x})\|_\infty^2 \geq \mu_1 (f(\mathbf{x}) - f^*)$$

Show that this new method achieves a linear convergence rate, i.e.

$$f(\mathbf{x}_t) - f^* \leq \left(1 - \frac{\mu_1}{L}\right)^t [f(\mathbf{x}_0) - f^*]$$



Solution: We have the following inequality from coordinate-wise Lipschitz-continuous gradient: +1/16/45+

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L} |\nabla_{i_t} f(\mathbf{x}_t)|^2$$

According to our sampling rule, we have $|\nabla_{i_t} f(\mathbf{x}_t)| = \|\nabla f(\mathbf{x})\|_\infty$, thus

$$\begin{aligned} f(\mathbf{x}_{t+1}) &\leq f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|_\infty^2 \\ &\leq f(\mathbf{x}_t) - \frac{\mu_1}{L} (f(\mathbf{x}_t) - f^*) \end{aligned}$$

Subtracting f^* from both sides we have:

$$f(\mathbf{x}_{t+1}) - f^* \leq \left(1 - \frac{\mu_1}{L}\right) (f(\mathbf{x}_t) - f^*)$$

By the equivalence between norms, we have that $\mu/d \leq \mu_1$, so this is faster than the rate with random selection.

General Understanding of the PL Inequality

Question 34: 3 points. Let A be some non-zero matrix. Define $f(\mathbf{x}) := g(A\mathbf{x})$ for a σ -strongly convex function g . Derive the constant for the PL inequality for f .

☐ 0 ☐ 1 ☐ 2 ☒ 3

Solution: Define $\mathbf{u} := A\mathbf{x}$ and $\mathbf{v} := A\mathbf{y}$. By the SC of g , we have

$$g(\mathbf{v}) \geq g(\mathbf{u}) + \nabla g(\mathbf{u})^\top (\mathbf{v} - \mathbf{u}) + \frac{\sigma}{2} \|\mathbf{v} - \mathbf{u}\|^2$$

Substitute $\mathbf{u} = A\mathbf{x}$ and $\mathbf{v} = A\mathbf{y}$ back, we have

$$g(A\mathbf{y}) \geq g(A\mathbf{x}) + \nabla g(A\mathbf{x})^\top (A\mathbf{y} - A\mathbf{x}) + \frac{\sigma}{2} \|A\mathbf{y} - A\mathbf{x}\|^2$$

Given $\nabla f(\mathbf{x}) = A^\top \nabla g(A\mathbf{x})$, we have

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{\sigma}{2} \|A(\mathbf{y} - \mathbf{x})\|^2$$

Using \mathbf{x}_p to denote the projection of \mathbf{x} onto the optimal solution set \mathcal{X}^* , we have

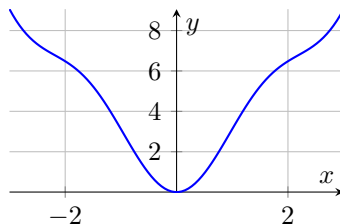
$$\begin{aligned} f(\mathbf{x}_p) &\geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{x}_p - \mathbf{x}) + \frac{\sigma}{2} \|A(\mathbf{x}_p - \mathbf{x})\|^2 \\ &\geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{x}_p - \mathbf{x}) + \frac{\sigma \theta(A)}{2} \|\mathbf{x}_p - \mathbf{x}\|^2 \\ &\geq f(\mathbf{x}) + \min_{\mathbf{y}} \left[\nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{\sigma \theta(A)}{2} \|\mathbf{y} - \mathbf{x}\|^2 \right] \\ &= f(\mathbf{x}) - \frac{1}{2\theta(A)\sigma} \|\nabla f(\mathbf{x})\|^2 \end{aligned}$$

where $\theta(A)$ is the smallest non-zero singular value of A .

Question 35: 2 points. If f satisfies the PL condition, does this imply that f is strongly convex? Justify your answer.

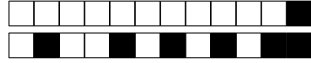
☐ 0 ☐ 1 ☒ 2

Solution: No, the PL condition does not imply convexity. For example, the function $f(x) = x^2 + 3\sin^2(x)$ satisfies the PL condition but is not convex. PL implies invexity though – all stationary points are global optimum.





DRAFT



DRAFT



DRAFT