



Exam Optimization for Machine Learning – CS-439
Prof. Martin Jaggi

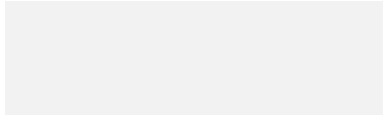


Fr. 6. July 2018 - 16h15 to 19h15, in CE1515

ID













STUDENT NAME

SCIPER: **SCIPER**

Signature: 

Wait for the start of the exam before turning to the next page. This document is printed double sided, 18 pages.

- 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 (handwritten or 11pt min font size) 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.

Newton-Raphson method

An easy method for computing the square root of a real number $y > 0$ by hand is as follows:

- (i) find x_0 , such that $x_0^2 \approx y$ (e.g. $y = 17$, $x_0 = 4$).
- (ii) Calculate the difference $d = y - x_0^2$ (e.g. $d = 17 - 4^2 = 1$).
- (iii) Output $x_1 = x_0 + \frac{d}{2x_0}$ (e.g. $x_1 = 4 + \frac{1}{8} = 4.125$).
- (iv) repeat (ii)–(iii) for higher accuracy.

This is an instance of the Newton-Raphson method, which defines the sequence $\{x_t\}_{t \geq 0}$ of real numbers by the following equation:

$$x_{t+1} := x_t - \frac{f(x_t)}{f'(x_t)}. \quad (1)$$

Question 1 What is the function $f(z)$ of which we aim to find a zero in the example above?

- $\sqrt{z - 17}$
 z^2
 $z^2 - 17$
 \sqrt{z}

Question 2 Now, suppose we are not happy with the solution $x_1 = 4.125$, because $x_1^2 = 17.015625$ is not accurate enough. What is the next iterate x_2 in the sequence (for $y = 17$, $x_0 = 4$ as above)? Use the following values: $\frac{17.015625}{4.125} = 4.125$, $\frac{0.015625}{4.125} \approx 0.0038$.

- 4.1288
 4.1269
 4.1212
 4.1231

Question 3 How many iterations do you (roughly) have to perform to compute the correct 16 significant digits in the above example ($y = 17$, $x_0 = 4$).

- 10^{16}
 $10^{\sqrt{16}} = 10^4$
 16
 $\frac{16}{2} = 8$
 $\sqrt{16} = 4$



Question 4 The Newton-Raphson method to find zeros of f can be interpreted as a second-order optimization method. Of course, one could also use the gradient method instead. How would the iterates of this scheme look like? (For carefully chosen stepsize γ).

- $x_{t+1} := x_t - \gamma f(x_t)$
- $x_{t+1} := x_t - \gamma f'(x_t)$
- $x_{t+1} := x_t - \gamma x_t$
- $x_{t+1} := x_t - \gamma (f'(x_t))^{-1} f(x_t)$
- $x_{t+1} := x_t - \gamma (f''(x_t))^{-1} f'(x_t)$

Newton's second-order optimization method

Question 5 As studied in the class, the update step for Newton's optimization method for an objective function $g : \mathbb{R}^n \rightarrow \mathbb{R}$ is given by

$$\mathbf{x}_{t+1} := \mathbf{x}_t - \nabla^2 g(\mathbf{x}_t)^{-1} \nabla g(\mathbf{x}_t)$$

For $n = 1$, how does this optimization method relate to the Newton-Raphson method from Equation (1) from the previous section?

- $f = g'$
- $f'' = g$
- $f' = g$
- $f = g''$

Question 6 Given a quadratic function $g : \mathbb{R}^n \rightarrow \mathbb{R}$ of the form $g(\mathbf{x}) = -\frac{1}{2} \mathbf{x}^\top A \mathbf{x} + \mathbf{b}^\top \mathbf{x} + c$ where $A \in \mathbb{R}^{n \times n}$ is a symmetric matrix. What are necessary and sufficient conditions for g to be convex?

- $-A$ positive semidefinite, and \mathbf{b} is non-negative
- The Hessian of g is negative definite for all \mathbf{x} , and \mathbf{b} is non-negative
- A positive semidefinite
- The Hessian of g is negative definite for all \mathbf{x}
- The Hessian of g is positive definite for all \mathbf{x}
- A positive semidefinite, and \mathbf{b} is non-negative
- $-A$ positive semidefinite
- The Hessian of g is positive definite for all \mathbf{x} , and \mathbf{b} is non-negative



Coordinate Descent

Question 7 Consider the least squares objective function

$$f(\mathbf{x}) := \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|^2, \quad (2)$$

for A a $m \times n$ matrix, $A = [\mathbf{a}_1, \dots, \mathbf{a}_n]$ with columns \mathbf{a}_i .

What is the gradient $\nabla f(\mathbf{x})$?

- A
- $A^\top \mathbf{Ax}$
- $A^\top (\mathbf{Ax} - \mathbf{b})$
- $A(A^\top \mathbf{x} - \mathbf{b})$

Question 8 We are now interested in the complexity of computing the gradient of f as in Equation (2). Each addition or multiplication of two real numbers counts as one operation. How expensive is it to compute the full gradient, given \mathbf{x} :

(Note here $\Theta(k)$ refers to a function growing at least and at most as fast as k in the variables of concern)

- $\Theta(n + m)$
- $\Theta(n^2 m^2)$
- $\Theta(n)$
- $\Theta(mn)$
- $\Theta(m^2 n)$
- $\Theta(mn^2)$
- $\Theta(m)$

Question 9 How expensive is it to compute just a single coordinate of the gradient of f as in Equation (2), given \mathbf{x} :

(Note here $\Theta(k)$ refers to a function growing at least and at most as fast as k in the variables of concern)

- $\Theta(n)$
- $\Theta(m^2 n)$
- $\Theta(m)$
- $\Theta(n^2 m^2)$
- $\Theta(n + m)$
- $\Theta(mn)$
- $\Theta(mn^2)$

Question 10 The complexity of Coordinate Descent depends on the coordinate-wise smoothness constants L_i . What is L_i for f as in Equation (2)?

- $L_i = \lambda_{\max}(A^\top A)$
- $L_i = \|\mathbf{a}_i\|^2$
- $L_i = \lambda_{\max}(A^\top A)/n$
- $L_i = \|A^\top A\|$
- $L_i = \|\mathbf{a}_i\|$



Frank-Wolfe

Consider the linear minimization oracle (LMO) for matrix completion, that is for

$$\min_{Y \in X \subseteq \mathbb{R}^{n \times m}} \sum_{(i,j) \in \Omega} (Z_{ij} - Y_{ij})^2$$

when $\Omega \subseteq [n] \times [m]$ is the set of observed entries from a given matrix Z . Our optimization domain X is the unit ball of the trace norm (or nuclear norm), which is known to be the convex hull of the rank-1 matrices

$$X := \text{conv}(\mathcal{A}) \quad \text{with} \quad \mathcal{A} := \left\{ \mathbf{u}\mathbf{v}^\top \mid \begin{array}{l} \mathbf{u} \in \mathbb{R}^n, \|\mathbf{u}\|_2=1 \\ \mathbf{v} \in \mathbb{R}^m, \|\mathbf{v}\|_2=1 \end{array} \right\}.$$

Question 11 Consider the LMO for this set X for a gradient at iterate $Y \in \mathbb{R}^{n \times m}$ (derive it if necessary). Compare the computational operation (or cost) needed to compute the LMO, as opposed to computing the *projection* onto X ?

Hint: Assume that the Singular Value Decomposition of a $n \times m$ matrix takes time $\Theta(n^2m)$, and computing the top singular vector takes time $\Theta(nm)$.

- LMO and projection both take $\Theta(n^2m)$
- LMO takes $\Theta(nm)$, and projection takes $\Theta(n^2m)$
- LMO takes $\Theta(n^2m)$, and projection takes $\Theta(nm)$

Smoothness and Strong Convexity

Consider an iterative optimization procedure.

Question 12 Which one of the following three inequalities is valid for a *smooth* convex function f :

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

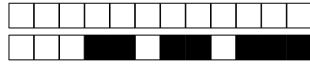
Question 13 Which one of the following three inequalities is valid for a *strongly convex* function f :

- $f(\mathbf{x}_t) - f(\mathbf{x}^*) \geq \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*) + \frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}^*\|^2$
- $f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*) - \frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}^*\|^2$
- $f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*) + \frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}^*\|^2$

Random search

Question 14 Consider derivative-free random search, with line-search, as discussed in the lecture.

- For strongly convex functions, random search converges as $\mathcal{O}(L \log(1/\varepsilon))$
- For convex functions, random search converges as $\mathcal{O}(dL/\varepsilon)$
- For convex functions, random search converges as $\mathcal{O}(dL \log(1/\varepsilon))$



Empirical comparison of different methods

Donald Duck's three nephews Huey, Dewey, and Louie have enrolled in CS 439. For their course project, they analyzed three different algorithms, namely *Gradient Descent*, *Accelerated Gradient Method* and *Newton's second-order optimization method* on a strongly convex optimization problem and plotted the performance of the algorithms on a graph. However, as it turns out they forgot to put a legend in their graph and due to some bug in their code, they plotted a line which corresponds to none of the algorithms. Can you help them in labelling their unlabelled graph?

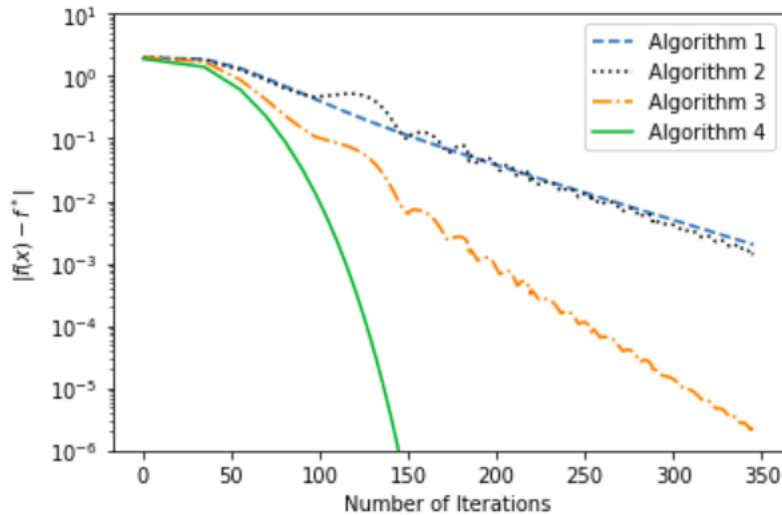


Figure 1: Performance of different optimization algorithms.

Question 15 Which optimization method corresponds to the error-curve for **Algorithm 1**?

- None
- Gradient Descent (with correct stepsize)
- Accelerated Gradient Method (with correct parameters)
- Newton's optimization method

Question 16 Which optimization method corresponds to the error-curve for **Algorithm 2**?

- None
- Newton's optimization method
- Accelerated Gradient Method (with correct parameters)
- Gradient Descent (with correct stepsize)

Question 17 Which optimization method corresponds to the error-curve for **Algorithm 3**?

- Accelerated Gradient Method (with correct parameters)
- Newton's optimization method
- None
- Gradient Descent (with correct stepsize)



Question 18 Which optimization method corresponds to the error-curve for **Algorithm 4**?

- Newton's optimization method
- None
- Accelerated Gradient Method (with correct parameters)
- Gradient Descent (with correct stepsize)

DRAFT



Second part, true/false questions

Question 19 (Convexity) The *epigraph* of a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is defined as

$$\text{epi}(f) := \{(\mathbf{x}, \alpha) \in \mathbb{R}^{d+1} \mid \mathbf{x} \in \text{dom}(f), \alpha \leq f(\mathbf{x})\},$$

TRUE FALSE

Question 20 (Convexity) The triangle inequality and homogeneity of a norm together imply that any norm is convex.

TRUE FALSE

Question 21 (Convex Sets) We consider C_1 and C_2 two convex sets in \mathbb{R}^d .

We define $C_1 + C_2 := \{x_1 + x_2, x_1 \in C_1, x_2 \in C_2\}$.

Is $C_1 + C_2$ a convex set?

TRUE FALSE

Question 22 (Differentiability) The function $(\max(0, x))^2$ is differentiable over \mathbb{R} .

TRUE FALSE

Question 23 (Coordinate Descent) Consider Coordinate Descent on a strongly convex and smooth objective function, depending on the coordinate-wise smoothness constants L_i (i.e. different Lipschitz constants for each gradient coordinate).

If we sample coordinate i with probability proportional to L_i , and use stepsize L_i , convergence is typically faster than uniform CD

TRUE FALSE

Question 24 (Coordinate Descent) In the same setting, if we sample coordinate i uniformly, and use stepsize $1/L_i$, convergence is typically faster than CD with fixed stepsize

TRUE FALSE



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.

Intersection of Convex Sets

We are given n convex sets $C_0 = \{C_1, \dots, C_n\}$ where each set $C_i \subseteq \mathbb{R}^d$. We want to design an algorithm which can check if the intersection of *all* of these sets is null i.e. we want to check if

$$\bigcap_{C_i \in C_0} C_i = \emptyset.$$

However we do not care about solving this problem exactly, but have a small leeway of magnitude $\varepsilon \geq 0$. To make this more mathematically precise, let us define some notation.

The distance between a set $C \subseteq \mathbb{R}^d$ and any point $\mathbf{y} \in \mathbb{R}^d$ is defined as

$$d(C, \mathbf{y}) \stackrel{\text{def}}{=} \min_{\mathbf{w} \in C} \|\mathbf{w} - \mathbf{y}\|_2.$$

We only want to distinguish between the following two cases for any $\varepsilon > 0$:

- (N) The intersection of the sets is non-empty, i.e. $\bigcap_{i=1}^n C_i \neq \emptyset$.
- (E) For any point $\mathbf{x} \in \mathbb{R}^d$, $\max_{i \in \{1, \dots, n\}} d(C_i, \mathbf{x}) \geq \varepsilon$.

We want to solve this problem using calls to an oracle which can compute the projection onto $C_i \in \mathcal{C}$. Let us define the projection oracle $P_i(\mathbf{x})$ for any $i \in \{1, \dots, n\}$ and $\mathbf{x} \in \mathbb{R}^d$ as

$$P_i(\mathbf{x}) := \operatorname{argmin}_{\mathbf{y} \in C_i} \|\mathbf{y} - \mathbf{x}\|_2.$$

We want to make as few calls to the projection oracle as possible. Our strategy will be to i) define a loss function and ii) run gradient descent. Then using our knowledge of convergence of gradient descent, we can argue about the number of oracle calls required.

First Approach.

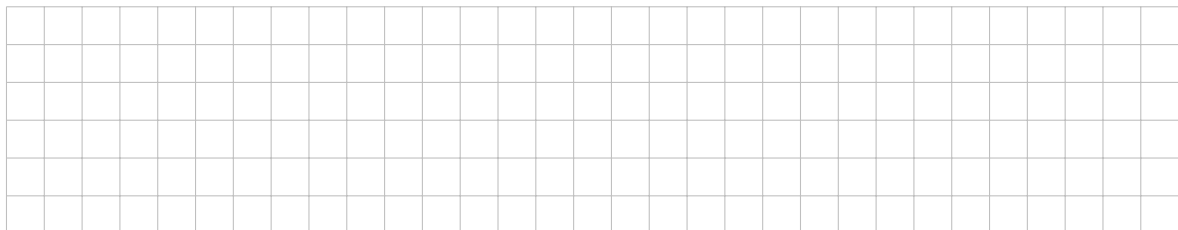
Inspired by the condition in case (E), let us define the following loss function:

$$g(\mathbf{x}) := \max_{i \in \{1, \dots, n\}} d(C_i, \mathbf{x}).$$

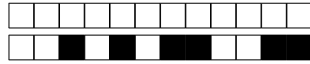
Question 24: 2 points. Is the function $g(\mathbf{x})$ convex? Is it Lipschitz?

Hint: maximum of convex functions is also convex.

0 1 2



Solution: The function $d(C_i, \mathbf{x})$ is convex and since maximum of convex functions is convex, $g(\mathbf{x})$ is also convex. Since for all i , $d(C_i, \mathbf{x})$ is 1-Lipschitz, the function $g(\mathbf{x})$ is also 1-Lipschitz.



Question 25: 5 points. What is the sub-gradient of g ? How many calls to the gradient oracle are needed to compute $\partial g(\mathbf{x})$ and $g(\mathbf{x})$?

Hint: Show that for two convex functions $g_1(x)$ and $g_2(x)$, $\partial g_i(x)$ is a subgradient in the set $\partial \max(g_1(x), g_2(x))$ where $g_i(x) := \max(g_1(x), g_2(x))$.

0 1 2 3 4 5

Solution: Since $d(C_i, \mathbf{x})$ can be written as $\|\mathbf{x} - P_i(\mathbf{x})\|$, the subgradient of $d(C_i, \mathbf{x})$ is

$$\partial d(C_i, \mathbf{x}) = (\mathbf{x} - P_i(\mathbf{x})) / \|\mathbf{x} - P_i(\mathbf{x})\|_2 .$$

Then the subgradient of $g(\mathbf{x})$ is

$$\partial g(\mathbf{x}) = \partial d(C_j, \mathbf{x}) = (\mathbf{x} - P_j(\mathbf{x})) / \|\mathbf{x} - P_j(\mathbf{x})\|_2 ,$$

where $j \in \{1, \dots, n\}$ is such that

$$j = \operatorname{argmax}_{i \in \{1, \dots, n\}} d(C_i, \mathbf{x}) .$$

Thus computing $\partial g(\mathbf{x})$ and $g(\mathbf{x})$ requires n calls to the projection oracle—same as that required in the case of $f(\mathbf{x})$.

Question 26: 5 points. Assume you are given a starting point \mathbf{x}_0 and a constant R such that $\|\mathbf{x}_0 - \mathbf{x}^*\|_2 \leq R$. Give the update step of gradient descent with an appropriate step-size. Show using the convergence of gradient descent we proved in class that for any optimum \mathbf{x}^* of g ,

$$\min_{t \in \{0, \dots, T\}} g(\mathbf{x}_t) - g(\mathbf{x}^*) \leq \frac{R}{\sqrt{T}} .$$

0 1 2 3 4 5

Solution: Since the norm of the gradient is less than 1, we will run gradient descent for T steps with step size $\gamma = \frac{R}{\sqrt{T}}$. The update is

$$\mathbf{x}_{t+1} \leq \mathbf{x}_t - \frac{R}{\sqrt{T}} \frac{\mathbf{x}_t - P_j(\mathbf{x}_t)}{\|\mathbf{x}_t - P_j(\mathbf{x}_t)\|_2} .$$

The convergence result for bounded gradients gives us that

$$\min_{t \in \{0, \dots, T\}} g(\mathbf{x}_t) - g(\mathbf{x}^*) \leq \frac{1}{T} \sum_{t=1}^T g(\mathbf{x}_t) - g(\mathbf{x}^*) \leq \frac{R}{\sqrt{T}} .$$

Question 27: 4 points. Using the result from the previous question, show that $\mathcal{O}(n/\varepsilon^2)$ calls to the projection oracle is sufficient to distinguish between case (N) and case (E) for our problem.

0 1 2 3 4

Solution: Run our algorithm for $T = R/\varepsilon^2$ steps. This requires Tn calls to the projection oracle. Then compute

$$\mathbf{z}_T = \operatorname{argmin}_{t \in \{0, \dots, T\}} g(\mathbf{x}_t) ,$$

using Tn additional calls to the projection oracle. If $g(\mathbf{z}_T) \leq \varepsilon$ then output case (N), else output case (E). This requires a total of $2nR/\varepsilon^2$.

The correctness follows because if case (E) were true then for any \mathbf{x} , $g(\mathbf{x}) > \varepsilon$ and if case (N) were true then $g(\mathbf{x}^*) = 0$ and we can use the result from Part 3.



Question 28: 6 points.

0 1 2 3 4 5 6

The convergence of the Frank-Wolfe algorithm was analyzed in class for only *smooth* functions. In this question we will examine if smoothness is necessary. Consider the following non-smooth function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$:

$$f(w, v) := \max \{w, v\} ,$$

restricted to a ball of radius 2 around the origin. We are then interested in finding

$$(w^*, v^*) := \operatorname{argmin}_{w^2+v^2 \leq 2} (\max \{w, v\}) .$$

Suppose we start at the origin $(0, 0)$ and run the Frank-Wolfe algorithm (with any step size rule). Since the function is not smooth, we will call the LMO oracle using an arbitrary subgradient instead of the gradient. Does this algorithm converge to the optimum?

Hint: First show that the iterates of Frank-Wolfe always lie in the convex hull of the starting point and the solutions of the LMO oracle.

Solution: It does not converge.

Let us first prove the hint by induction over the iterations t . Suppose that at iteration t , \mathbf{x}_t is the current iterate and \mathbf{s}_t is the output of the LMO oracle. For $t = 0$, the statement is trivially true. Now for any $t \geq 0$ by induction assume that \mathbf{x}_t is a convex combination of \mathbf{x}_0 and $\{\mathbf{s}_0, \dots, \mathbf{s}_{t-1}\}$. The next iterate \mathbf{x}_{t+1} is computed as

$$\mathbf{x}_{t+1} = (1 - \gamma_t)\mathbf{x}_t + \gamma_t\mathbf{s}_t ,$$

for some step-size $\gamma_t \in [0, 1]$. Thus \mathbf{x}_{t+1} is a convex combination of \mathbf{x}_0 and $\{\mathbf{s}_0, \dots, \mathbf{s}_t\}$.

Let us compute the solution of the LMO for our problem. At any point, either $(1, 0)$ or $(0, 1)$ is a subgradient of $f(w, v)$. Using the LMO with each of these subgradients gives

$$\operatorname{argmin}_{w^2+v^2 \leq 2} \{(1, 0) \cdot (w, v)^\top\} = (-\sqrt{2}, 0) ,$$

and

$$\operatorname{argmin}_{w^2+v^2 \leq 2} \{(0, 1) \cdot (w, v)^\top\} = (0, -\sqrt{2}) .$$

This means that at every step, the iterates of Frank-Wolfe will always lie in the convex hull of $0, 0$, $(-\sqrt{2}, 0)$, and $(0, \sqrt{2})$. The optimum is $(-1, -1)$ and is outside this convex hull. Thus the algorithm will never converge to the optimum even if run for a very long time. \square



Newton's second-order optimization method

As studied in the class, the update step for Newton's optimization method for an objective function $g : \mathbb{R}^n \rightarrow \mathbb{R}$ is given by

$$\mathbf{x}_{t+1} := \mathbf{x}_t - \nabla^2 g(\mathbf{x}_t)^{-1} \nabla g(\mathbf{x}_t)$$

Question 29: *2 points.* What happens when Newton's optimization method is run on a convex quadratic function? Explain.

0 1 2



Question 30: *2 points.* **Affine Invariance of the Newton's method**

Consider $h(\mathbf{x}) := g(M\mathbf{x})$ where $M \in \mathbb{R}^{n \times n}$ is invertible where g is some convex function. Show that the Newton steps for h and g are also related by the same linear transformation, i.e., $\Delta \mathbf{x}_t = M \Delta \mathbf{y}_t$ where $\Delta \mathbf{x}_t$ and $\Delta \mathbf{y}_t$ are the Newton steps at the t^{th} iteration for h and g respectively. We assume $\mathbf{x}_0 = M\mathbf{y}_0$ are the starting iterates for h and g respectively.

0 1 2





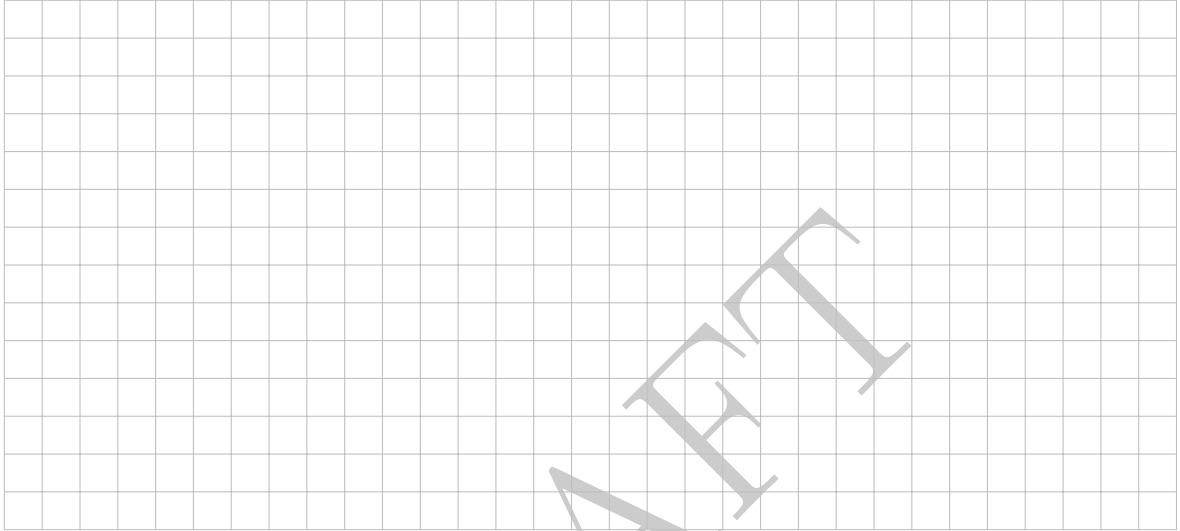
Coordinate Descent

Question 31: 2 points. Given a matrix A , we define $\lambda_{\min}(A^T A)$ and $\lambda_{\max}(A^T A)$ to be the smallest and largest eigenvalues of $A^T A$.

Show that for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$,

$$\lambda_{\min}(A^T A) \|\mathbf{x} - \mathbf{y}\|^2 \leq \|A(\mathbf{x} - \mathbf{y})\|^2 \leq \lambda_{\max}(A^T A) \|\mathbf{x} - \mathbf{y}\|^2.$$

0 1 2

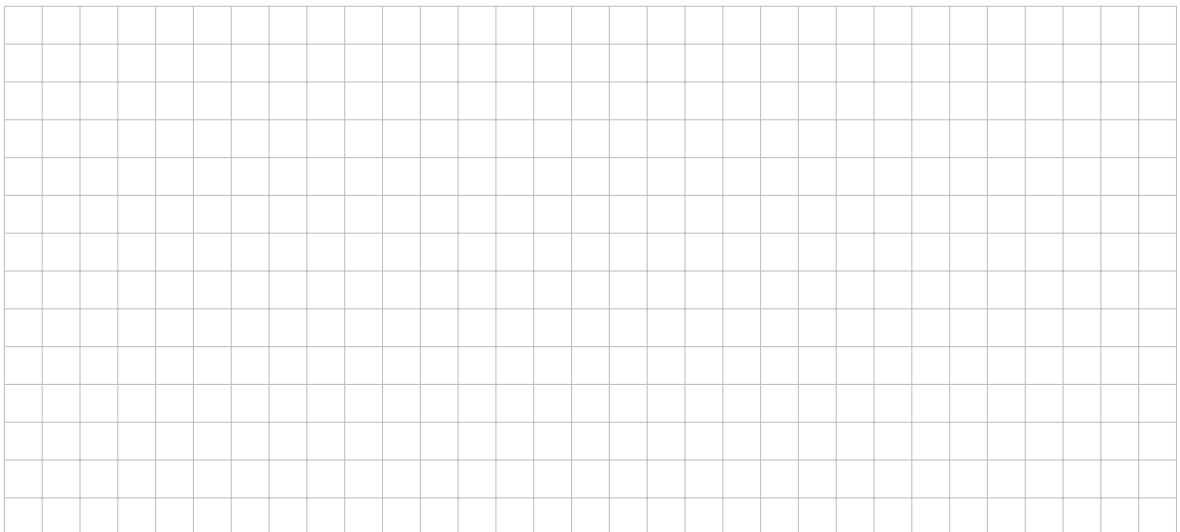


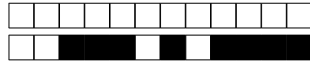
Question 32: 3 points. Show that for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, for any $\mathbf{b} \in \mathbb{R}^n$:

$$\|A\mathbf{x} - \mathbf{b}\|^2 \leq \|A\mathbf{y} - \mathbf{b}\|^2 + [A^T(A\mathbf{y} - \mathbf{b})]^T(\mathbf{x} - \mathbf{y}) + \frac{\lambda_{\max}(A^T A)}{2} \|\mathbf{x} - \mathbf{y}\|^2.$$

What does that imply for $f(\mathbf{x}) := \|A\mathbf{x} - \mathbf{b}\|^2$?

0 1 2 3



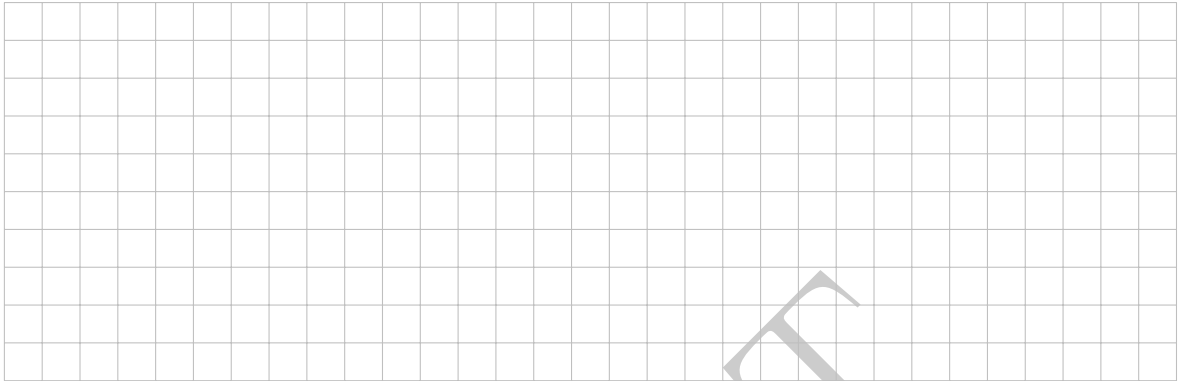


Question 33: 2 points. For $f(\mathbf{x}) := \|\mathbf{Ax} - \mathbf{b}\|^2$, we now perform one step of coordinate descent. I.e. for a given point $\mathbf{x}_t \in \mathbb{R}^n$ we do a step of the form

$$\mathbf{x}_{t+1} := \mathbf{x}_t - \gamma_t (\nabla f(\mathbf{x}_t))_i \cdot \mathbf{e}_i$$

where $\mathbf{e}_i \in \mathbb{R}^n$ denotes a standard unit vector. For i fixed, compute the best γ_t .

0 1 2



Smooth strongly convex SGD

We consider a function $f(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x})$ on \mathbb{R}^d , and we assume that the functions f_i are convex, differentiable.

We furthermore assume that f is L -smooth, that is that ∇f is L -Lipschitz.

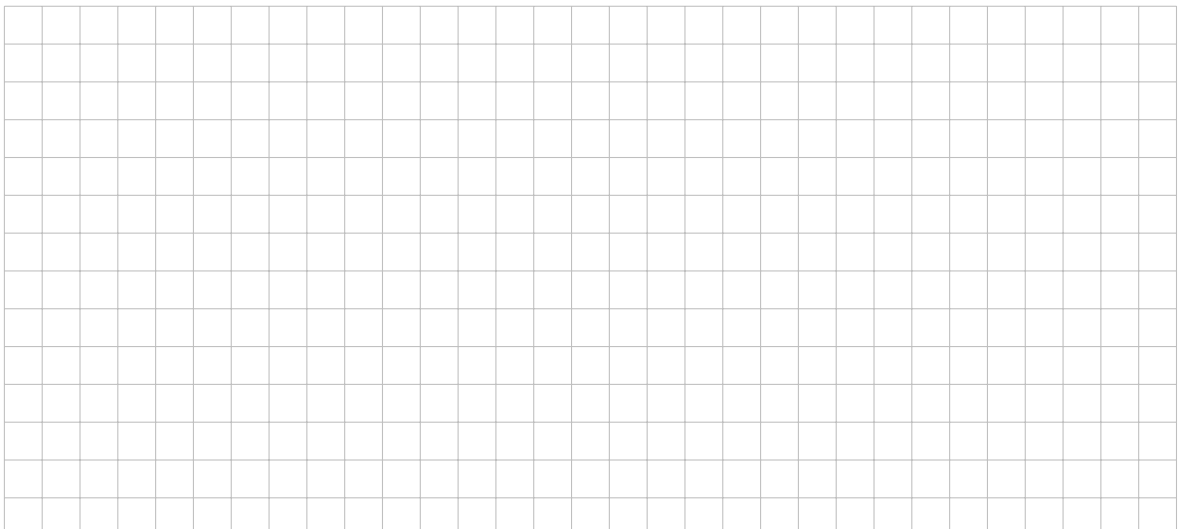
We consider SGD defined as the following algorithm: Let $\mathbf{x}_0 \in \mathbb{R}^d$, and for any $t \geq 1$, for a sequence of step sizes γ_t , define

$$\mathbf{x}_{t+1} := \mathbf{x}_t - \gamma_t \mathbf{g}_t.$$

We first consider $\mathbf{g}_t := \nabla f_{i_t}(\mathbf{x}_t)$, with i_t uniformly and independently sampled from $\{1, \dots, n\}$.

Question 34: 2 points. Show that \mathbf{g}_t is an unbiased estimator of the gradient $\nabla f(\mathbf{x}_t)$.

0 1 2





Question 35: 6 points. Combining the two valid equations of smoothness and strong convexity (as also stated in Questions 12 and 13), prove in detailed steps that, if $\gamma_t \leq \frac{1}{L}$, SGD in this setting converges as

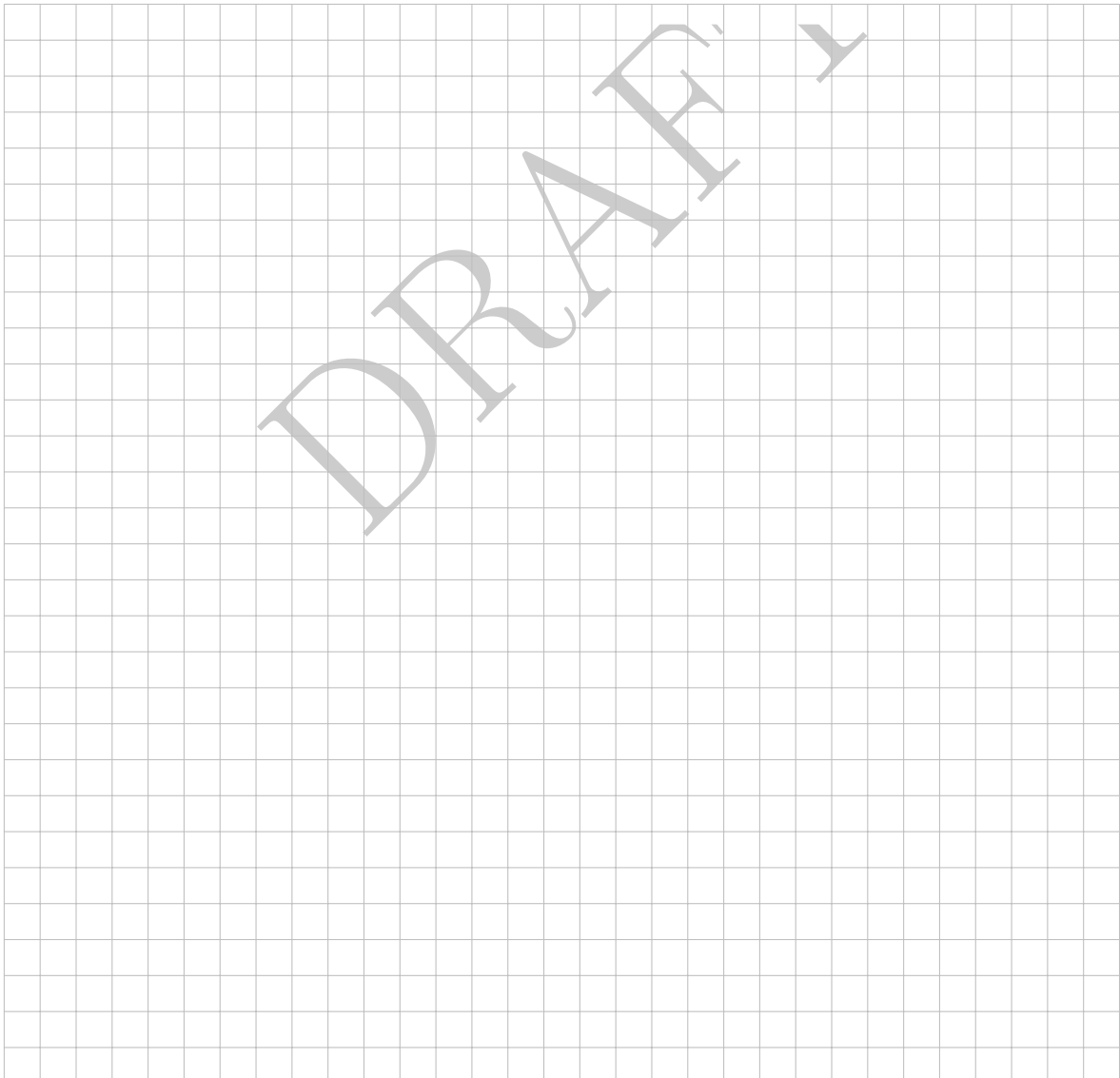
$$\mathbb{E} [f(\mathbf{x}_{t+1}) - f(\mathbf{x}^*)] \leq \gamma_t \mathbb{E} [\|\mathbf{g}_t - \nabla f(\mathbf{x}_t)\|^2] + \frac{(1 - \gamma_t \mu) \mathbb{E} [\|\mathbf{x}_t - \mathbf{x}^*\|^2] - \mathbb{E} [\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2]}{2\gamma_t}. \tag{3}$$

For comparison, recall the following result from Lecture 6 (slide 6):

$$\mathbb{E} [f(\mathbf{x}_{t+1}) - f(\mathbf{x}^*)] \leq \frac{\gamma_t B^2}{2} + \frac{(1 - \gamma_t \mu) \mathbb{E} [\|\mathbf{x}_t - \mathbf{x}^*\|^2] - \mathbb{E} [\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2]}{2\gamma_t}. \tag{4}$$

under the bounded gradient assumption $\mathbb{E}[\|\mathbf{g}_t\|^2] \leq B^2$.

How do the two results compare?



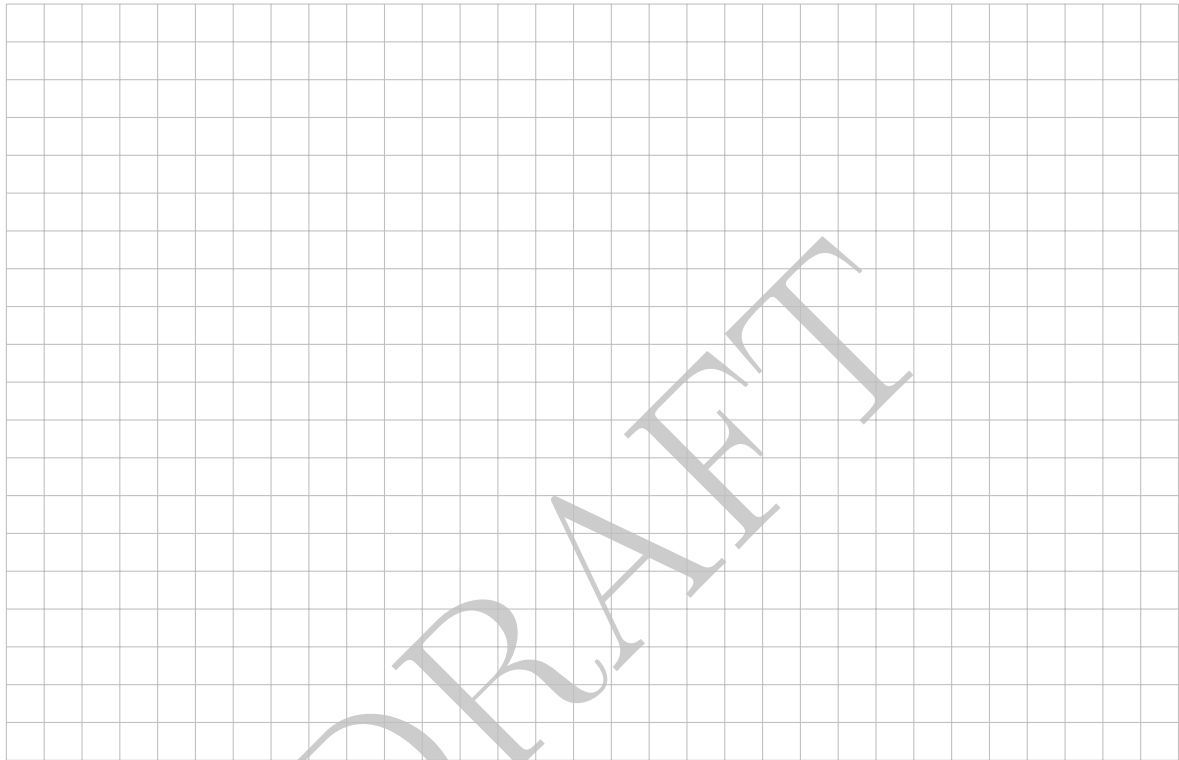


Question 36: 4 points. Recall the possible choices of learning rate (γ_t) in the situation of the previous question. What is the resulting rate of convergence? Which estimator do we eventually consider?

Comment on the assumption $\gamma_t \leq \frac{1}{L}$. Is it a restriction? Which choice of step size could be used,

- a) for getting $\mathcal{O}(\log(t)/t)$ convergence, and
- b) for getting $\mathcal{O}(1/t)$?

0 1 2 3 4



Solution: $(\mu t)^{-1}$ to get $\log(t)/t$ we consider averaged estimator. Or with the trick seen in the class $\frac{1}{\mu(t+1)}$ with non uniform averaging, rate B^2/t . Here the bound on the step size does not hold for the first iterations. Considering $\gamma_t = \frac{1}{\alpha + \mu t}$ would be a solution (e.g. $\alpha = L$).