

**Instructions:**

- All your solutions should be prepared in L<sup>A</sup>T<sub>E</sub>X and the PDF and .tex should be submitted to Canvas. Please submit all your files as ONE archive of filetype zip, tgz, or tar.gz.
- Name the file [your-first-name]\_[your-last-name].[filetype]. For example, I would call my submission rasmus\_kyng.zip.
- INCLUDE your name in the submission pdf and any files with code.
- If the TFs cannot easily deduce your identity from your files alone, they may decide not to grade your submission.
- For each question, a well-written and correct answer will be selected a sample solution for the entire class to enjoy. If you prefer that we do not use your solutions, please indicate this clearly on the first page of your assignment.
- For this homework, you will need to understand the section material from section on Mar. 2nd. You can find the notes here: [http://rasmuskyng.com/am221\\_spring18/sections/sec6.pdf](http://rasmuskyng.com/am221_spring18/sections/sec6.pdf).

**1. Subgradients.** In this problem we consider a continuous convex function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ .

- Show that the subdifferential  $\partial f(\mathbf{x})$  of  $f$  at  $\mathbf{x} \in \mathbb{R}^n$  is a closed convex set of  $\mathbb{R}^n$ .
- Show that  $\mathbf{x}^* \in \mathbb{R}^n$  is a minimizer of  $f$  (i.e a solution to  $\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$ ) iff  $0 \in \partial f(\mathbf{x})$ .

a.

$$\begin{aligned}
 \partial f(\mathbf{x}) &= \{g : g \text{ is a subgradient of } f \text{ at } \mathbf{x}\} \\
 &= \{g : f(\mathbf{y}) \geq f(\mathbf{x}) + g^\top(\mathbf{y} - \mathbf{x}), \forall \mathbf{y}\} \\
 &= \bigcap_{\mathbf{y} \in \mathbb{R}^n} \{g : f(\mathbf{y}) \geq f(\mathbf{x}) + g^\top(\mathbf{y} - \mathbf{x})\}
 \end{aligned}$$

Thus,  $\partial f(\mathbf{x})$  is an intersection of halfspaces. Hence,  $\partial f(\mathbf{x})$  is convex and closed.

- Let (i)  $\mathbf{x}^*$  is a minimizer of  $f$ , (ii)  $0 \in \partial f(\mathbf{x})$ .

**Proof of (i)  $\Rightarrow$  (ii)**

$$\begin{aligned}(i) &\Rightarrow f(\mathbf{y}) \geq f(\mathbf{x}^*), \forall \mathbf{y} \in \mathbb{R}^n \\ &\Rightarrow f(\mathbf{y}) \geq f(\mathbf{x}^*) + 0^\top(\mathbf{y} - \mathbf{x}^*) \\ &\Rightarrow 0 \in \partial f(\mathbf{x})\end{aligned}$$

**Proof of (ii)  $\Rightarrow$  (i)**

$$\begin{aligned}(ii) &\Rightarrow f(\mathbf{y}) \geq f(\mathbf{x}^*) \\ &\Rightarrow \mathbf{x}^* \text{ is a minimizer of } f\end{aligned}$$

**2. Perceptron revisited.** In this problem we will revisit the perceptron algorithm of Lecture 2 to find a separating hyperplane for a linearly separable dataset. The dataset  $\mathcal{D}$  is a set of pairs  $\mathcal{D} \stackrel{\text{def}}{=} \{(\mathbf{x}_i, y_i), 1 \leq i \leq n\}$  with  $\mathbf{x}_i \in \mathbb{R}^{d-1}$  and  $y_i \in \{0, 1\}$ . We saw that after transformation of the data, finding a separating hyperplane amounts to finding  $\mathbf{w} \in \mathbb{R}^d$  such that  $\mathbf{w}^\top \mathbf{x}'_i > 0$  for all  $i$ , where the definition of  $\mathbf{x}'_i$  using  $\mathbf{x}_i$  and  $y_i$  is given in the lecture notes.

Let us define the following function:

$$f(\mathbf{w}) \stackrel{\text{def}}{=} \sum_{i=1}^n \max(0, -\mathbf{w}^\top \mathbf{x}'_i), \mathbf{w} \in \mathbb{R}^d$$

- a. Show that  $f$  is convex and non-negative over  $\mathbb{R}^d$ . Is it differentiable?
- b. Assume that the dataset  $\mathcal{D}$  is linearly separable. Show that any  $\mathbf{w}^* \in \mathbb{R}^d$  solution to:

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$$

defines a separating hyperplane of  $\mathcal{D}$ . What is the value of  $f$  at  $\mathbf{w}^*$ ?

- c. Let us define:

$$f_i(\mathbf{w}) \stackrel{\text{def}}{=} \max(0, -\mathbf{w}^\top \mathbf{x}'_i), \mathbf{w} \in \mathbb{R}^d$$

and:

$$g_i(\mathbf{w}) = \begin{cases} 0 & \text{if } \mathbf{w}^\top \mathbf{x}'_i > 0 \\ -\mathbf{x}'_i & \text{otherwise} \end{cases}$$

show that for all  $\mathbf{w} \in \mathbb{R}^d$ ,  $g_i(\mathbf{w})$  is a subgradient of  $f_i$  at  $\mathbf{w}$ .

- d. Describe the perceptron algorithm in the language of subgradients and the algorithms for convex optimization we saw in class. In particular, how would you describe the normalization by  $\|\mathbf{x}'_i\|$  in the perceptron algorithm. Also note that each iteration of the perceptron focuses on one data point of the dataset at a time, can you draw an analogy with something we saw in class?

a.

$$f(\lambda \mathbf{w}_1 + (1 - \lambda) \mathbf{w}_2) = \sum_{i=1}^n \max(0, -(\lambda \mathbf{w}_1 + (1 - \lambda) \mathbf{w}_2)^\top \mathbf{x}'_i)$$

On the other hand,

$$\lambda f(\mathbf{w}_1) + (1 - \lambda) f(\mathbf{w}_2) = \sum_{i=1}^n \lambda \max(0, -\mathbf{w}_1^\top \mathbf{x}'_i) + (1 - \lambda) \max(0, -\mathbf{w}_2^\top \mathbf{x}'_i)$$

Now, let

$$f_i(\mathbf{w}) = \max(0, -\mathbf{w}^\top \mathbf{x}'_i)$$

$f_i(\mathbf{w})$  is convex because 0 and  $-\mathbf{w}^\top \mathbf{x}'_i$  are convex and  $\max$  is taking the intersection of two convex sets. Hence,

$$\max(0, -(\lambda \mathbf{w}_1 + (1 - \lambda) \mathbf{w}_2)^\top \mathbf{x}'_i) \geq \lambda \max(0, -\mathbf{w}_1^\top \mathbf{x}'_i) + (1 - \lambda) \max(0, -\mathbf{w}_2^\top \mathbf{x}'_i), \forall i = 1, \dots, n$$

Thus,

$$f(\lambda \mathbf{w}_1 + (1 - \lambda) \mathbf{w}_2) \geq \lambda f(\mathbf{w}_1) + (1 - \lambda) f(\mathbf{w}_2)$$

Therefore, have that  $f$  is convex.

$f$  is nonnegative because  $f_i(\mathbf{w}) \geq 0, \forall i$  and  $f$  is simply summing  $f_i$ 's.

$f$  is not differentiable ( $f$  is differentiable everywhere except 0).

b.

$$\mathbf{w}^* = \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{x}) \Rightarrow f(\mathbf{w}) \geq f(\mathbf{w}^*), \forall \mathbf{w}^* \in \mathbb{R}^d$$

Since

$$f(w) = \sum_{i=1}^n \max(0, -\mathbf{w}^\top \mathbf{x}'_i) \geq \sum_{i=1}^n 0 = 0$$

$f(\mathbf{w}^*) = 0$ . This is achieved when  $\max(0, -\mathbf{w}^\top \mathbf{x}'_i) = 0, \forall i$ , which implies  $\mathbf{w}^\top \mathbf{x}'_i > 0$ . Hence,  $\mathbf{w}^*$  defines a hyperplane of  $D$ .

c. If  $\mathbf{w}^\top \mathbf{x}'_i > 0$ ,

$$f_i(\mathbf{w}') = \max(0, -\mathbf{w}'^\top \mathbf{x}'_i)$$

$$f_i(\mathbf{w}) + g_i(\mathbf{w})^\top (\mathbf{w}' - \mathbf{w}) = 0$$

Hence,  $f_i(\mathbf{w}') \geq f_i(\mathbf{w}) + g_i(\mathbf{w})^\top (\mathbf{w}' - \mathbf{w})$ .

Else if  $\mathbf{w}^\top \mathbf{x}'_i \leq 0$ ,

$$f_i(\mathbf{w}') = \max(0, -\mathbf{w}'^\top \mathbf{x}'_i)$$

$$f_i(\mathbf{w}) + g_i(\mathbf{w})^\top (\mathbf{w}' - \mathbf{w}) = -\mathbf{w}'^\top \mathbf{x}'_i$$

Hence,  $f_i(\mathbf{w}') \geq f_i(\mathbf{w}) + g_i(\mathbf{w})^\top (\mathbf{w}' - \mathbf{w})$ .

So, either way, we have that  $g_i(\mathbf{w})$  is a subgradient of  $f_i$  at  $\mathbf{w}$ .

d. The perceptron algorithm is performing subgradient descent, instead of the gradients. The normalization by  $\|\mathbf{x}'_i\|$  can be thought of as defining a step size of the descent. The perceptron optimization is analogous to steepest descent we saw in class in that it is taking one data point in the dataset at a time.

**3. Gradient descent with weaker assumptions.** In this problem we will analyze the gradient descent algorithm of Lecture 9 under weaker regularity assumptions. We consider a differentiable convex function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  and only assume that  $f$ 's gradient is  $L$ -Lipschitz continuous as introduced in the previous problem set:

$$\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\| \leq L\|\mathbf{y} - \mathbf{x}\|, (\mathbf{x}, \mathbf{y}) \in \mathbb{R}^n \times \mathbb{R}^n$$

We do **not** assume that  $f$  is twice-differentiable. Finally, we choose a constant step size  $t = \frac{1}{L}$  instead of doing exact or backtracking line search.

- a. Show that the  $L$ -Lipschitz continuous assumption on  $f$ 's gradient implies the following quadratic upper bound on  $f$ :

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2, \mathbf{x} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^n$$

**Hint:** use the fact you proved in Problem set 5 that  $\frac{L}{2} \mathbf{x}^\top \mathbf{x} - f(\mathbf{x})$  is a convex function.

- b. Let us denote by  $\mathbf{x}^{(k)}$  the current solution at the  $k$ th iteration of the gradient descent algorithm. Show that:

$$f(\mathbf{x}^{(k+1)}) \leq f(\mathbf{x}^{(k)}) - \frac{1}{2L} \|\nabla f(\mathbf{x}^{(k)})\|^2, k \in \mathbb{N}$$

Show that this implies:

$$f(\mathbf{x}^{(k+1)}) \leq f(\mathbf{x}^*) + \nabla f(\mathbf{x}^{(k)})^\top (\mathbf{x}^{(k)} - \mathbf{x}^*) - \frac{1}{2L} \|\nabla f(\mathbf{x}^{(k)})\|^2, k \in \mathbb{N}$$

where  $\mathbf{x}^*$  is a minimizer of  $f$  over  $\mathbb{R}^n$ .

- c. Show that:

$$\nabla f(\mathbf{x}^{(k)})^\top (\mathbf{x}^{(k)} - \mathbf{x}^*) - \frac{1}{2L} \|\nabla f(\mathbf{x}^{(k)})\|^2 = \frac{L}{2} (\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2 - \|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|^2), k \in \mathbb{N}$$

hence, using b., we have:

$$f(\mathbf{x}^{(k+1)}) - f(\mathbf{x}^*) \leq \frac{L}{2} (\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2 - \|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|^2), k \in \mathbb{N}$$

- d. Show that part c. implies:

$$f(\mathbf{x}^{(k)}) - f(\mathbf{x}^*) \leq \frac{L}{2k} \|\mathbf{x}^0 - \mathbf{x}^*\|^2, k \in \mathbb{N}$$

Given  $\varepsilon > 0$ , how many iterations of gradient descent are required to obtain  $f(\mathbf{x}^{(k)}) - f(\mathbf{x}^*) < \varepsilon$ ? How does this compare to the strongly convex case?

a. Since  $g(\mathbf{x}) = \frac{L}{2}\mathbf{x}^\top\mathbf{x} - f(\mathbf{x})$  is convex,

$$\begin{aligned}
g(\mathbf{y}) &\geq g(\mathbf{x}) + \nabla g(\mathbf{x})^\top(\mathbf{y} - \mathbf{x}) \\
&\Leftrightarrow \frac{L}{2}\mathbf{x}^\top\mathbf{y} - f(\mathbf{y}) \leq \frac{L}{2}\mathbf{x}^\top\mathbf{x} - f(\mathbf{x}) + (L\mathbf{x} - \nabla f(\mathbf{x}))^\top(\mathbf{y} - \mathbf{x}) \\
&\Leftrightarrow f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top(\mathbf{y} - \mathbf{x}) + \frac{L}{2}(\mathbf{y}^\top\mathbf{y} - 2L\mathbf{x}^\top\mathbf{y} - \mathbf{x}^\top\mathbf{x}) \\
&\Leftrightarrow f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top(\mathbf{y} - \mathbf{x}) + \frac{L}{2}(\mathbf{y}^\top\mathbf{y} - 2L\mathbf{x}^\top\mathbf{y} + \mathbf{x}^\top\mathbf{x}) \\
&\Leftrightarrow f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top(\mathbf{y} - \mathbf{x}) + \frac{L}{2}\|\mathbf{y} - \mathbf{x}\|^2
\end{aligned}$$

b. Since  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \frac{1}{L}\nabla f(\mathbf{x}^{(k)})$ ,

$$\begin{aligned}
f(\mathbf{x}^{(k+1)}) &= f(\mathbf{x}^{(k)} - \frac{1}{L}\nabla f(\mathbf{x}^{(k)})) \\
&\leq f(\mathbf{x}^{(k)}) + \nabla f(\mathbf{x}^{(k)})^\top(-\frac{1}{L}\nabla f(\mathbf{x}^{(k)})) + \frac{L}{2}\|-\frac{1}{L}\nabla f(\mathbf{x}^{(k)})\|^2 \\
&= f(\mathbf{x}^{(k)}) - \frac{1}{2L}\|\nabla f(\mathbf{x}^{(k)})\|^2
\end{aligned}$$

Now, since  $f$  is convex,

$$\begin{aligned}
f(\mathbf{x}^*) &\geq f(\mathbf{x}^{(k)}) + \nabla f(\mathbf{x}^{(k)})^\top(\mathbf{x}^* - \mathbf{x}^{(k)}) \\
&\Leftrightarrow f(\mathbf{x}^{(k)}) \leq f(\mathbf{x}^*) + \nabla f(\mathbf{x}^{(k)})^\top(\mathbf{x}^{(k)} - \mathbf{x}^*)
\end{aligned}$$

Plugging this in to the derived inequality above,

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

c.

$$\begin{aligned}
&\nabla f(\mathbf{x}^{(k)})^\top(\mathbf{x}^{(k)} - \mathbf{x}^*) - \frac{1}{2L}\|\nabla f(\mathbf{x}^{(k)})\|^2 \\
&= \frac{L}{2}(\frac{2}{L}\nabla f(\mathbf{x}^{(k)})^\top(\mathbf{x}^{(k)} - \mathbf{x}^*) - \frac{1}{L^2}\|\nabla f(\mathbf{x}^{(k)})\|^2) \\
&= \frac{L}{2}\frac{1}{L}\nabla f(\mathbf{x}^{(k)})^\top(2(\mathbf{x}^{(k)} - \mathbf{x}^*) - \frac{1}{L}\nabla f(\mathbf{x}^{(k)})) \\
&= \frac{L}{2}((\mathbf{x}^{(k)} - \mathbf{x}^*) + (\mathbf{x}^{(k)} - \mathbf{x}^* - \frac{1}{L}\nabla f(\mathbf{x}^{(k)}))((\mathbf{x}^{(k)} - \mathbf{x}^*) - (\mathbf{x}^{(k)} - \mathbf{x}^* - \frac{1}{L}\nabla f(\mathbf{x}^{(k)}))) \\
&= \frac{L}{2}(\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2 - \|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|^2)
\end{aligned}$$

- d. Since  $f$  is convex,  $f(\mathbf{x}^{(k)})$  is nonincreasing for increasing steps,

$$\begin{aligned}
 f(\mathbf{x}^{(k)}) - f(\mathbf{x}^*) &\leq \frac{1}{k} \sum_{i=1}^k (f(\mathbf{x}^{(i)}) - f(\mathbf{x}^*)) \\
 &\leq \frac{L}{2k} \left( \sum_{i=1}^k \|\mathbf{x}^{(i-1)} - \mathbf{x}^*\|^2 - \|\mathbf{x}^{(i)} - \mathbf{x}^*\|^2 \right) \\
 &\leq \frac{L}{2k} (\|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2 - \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2) \\
 &\leq \frac{L}{2k} \|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2
 \end{aligned}$$

Hence, for  $f(\mathbf{x}^{(k)}) - f(\mathbf{x}^*) < \varepsilon$ , we need

$$k > \frac{L}{2\varepsilon} \|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2$$

**4. Gradient descent, condition number, Newton's method.** In this problem we will consider the following minimizing problem:

$$\min_{\mathbf{x} \in \mathbb{R}^2} f(\mathbf{x}) \stackrel{\text{def}}{=} \min_{\mathbf{x} \in \mathbb{R}^2} \mathbf{x}^\top A \mathbf{x}$$

with:

$$A = \begin{pmatrix} 1 + \lambda & 1 - \lambda \\ 1 - \lambda & 1 + \lambda \end{pmatrix}$$

where  $\lambda$  is a real number with  $\lambda \geq 1$ .

- Compute the gradient of  $f$ , its Hessian and its eigenvalues as a function of  $\lambda$ . What is the optimal solution of the above problem?
- Implement the gradient descent algorithm for the above problem. Use backtracking line search as seen in section with  $\alpha = 0.3$  and  $\beta = 0.7$ .
- Run the gradient descent algorithm for several values of  $\lambda$  between 1 and  $10^5$ . For each value of  $\lambda$ , record the number of iterations required to reach a solution with error smaller than  $10^{-10}$ . Choose  $x^0 = [1., 2.]$  as your starting point. Draw a plot of the number of iterations as a function of  $\lambda$ . How would you explain these results?
- Implement Newton's method for the above problem. Use backtracking line search with  $\alpha = 0.3$  and  $\beta = 0.7$ . For the same values of  $\lambda$  you used in c., show the number of the iterations required to reach the same error  $10^{-10}$ . How would you explain those results?

a.

$$\begin{aligned}
 \nabla f(\mathbf{x}) &= \frac{\partial \mathbf{x}^\top A \mathbf{x}}{\partial \mathbf{x}} \\
 &= (A + A^\top) \mathbf{x} \quad (\because \text{hw1}) \\
 &= 2A \mathbf{x}
 \end{aligned}$$

$$\begin{aligned} H_f(\mathbf{x}) &= \frac{2A\mathbf{x}}{\partial\mathbf{x}} \\ &= 2A \end{aligned}$$

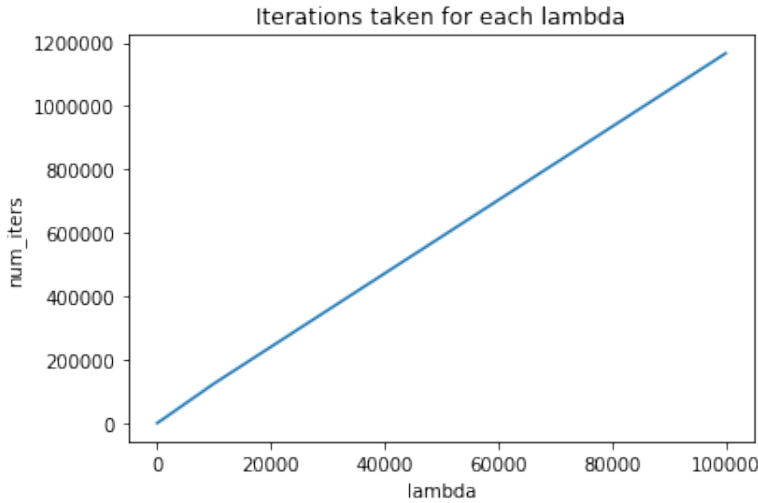
$$\begin{aligned} \det|A - \alpha I| &= (1 + \lambda - \alpha)^2 - (1 - \lambda)^2 \\ &= \alpha^2 - 2\alpha(1 + \lambda) + 4\lambda \end{aligned}$$

Hence, when  $\det|A - \alpha I| = 0 \Leftrightarrow \alpha = 1 + 2\lambda, 2$ . Since the eigen values are both positive, Hessian is positive definite. Hence, the optimization problem is strictly convex. So, the optimal solution is

$$\nabla f(\mathbf{x}) = 0 \Leftrightarrow \mathbf{x} = 0$$

b. See hw6.py for the implementation.

c. I ran gradient descent for  $[10^0, 10^1, 10^2, 10^3, 10^4, 10^5, 10^6]$



The result shows that as  $\lambda$  increases the number of iterations it takes the algorithm to converge increases (perhaps linearly). This is because when we have large  $\lambda$ , and starting from  $x^{(0)} = [1, 2]$ , which is an imbalanced start point, (1) the gradient is larger, making the algorithm to take longer time to approach 0 and (2) the learning rate  $t$  obtained by backline search is smaller, making the decrease in the gradient even slower every epoch.

d. I ran gradient descent for  $[10^0, 10^1, 10^2, 10^3, 10^4, 10^5]$  (because of run time issue).

Same as c, the result shows that as  $\lambda$  increases the number of iterations it takes the algorithm to converge increases (perhaps linearly). The reason is the same as c, but the effect is exacerbated by the fact that in the update, we're applying the inverse hessian to the update rule as well.

