

ECE 6270, Spring 2021

Homework #3

Due Sunday, Feb 14, at 11:59pm Suggested Reading: B&V, Sections 4.2.1-4.2.3 and 9.2-9.3 (You can skip the convergence analysis in 9.3 for now.)

1. Prepare a one paragraph summary of what we talked about in the last week of class. I do not want just a bulleted list of topics, I want you to use complete sentences and establish context (Why is what we have learned relevant? How does it connect with other classes?). The more insight you give, the better.
2. Provide feedback to your peers on Homework #2 in Canvas.
3. **Monotone gradients.** In class we showed that or a differentiable function f being convex is equivalent to the statement that

$$f(\mathbf{x}) \geq f(\mathbf{y}) + \langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{y}) \rangle, \quad (1)$$

for all $\mathbf{x}, \mathbf{y} \in \text{dom } f$. Here we will provide another equivalent characterization of convexity for differentiable f . Specifically, the first-order condition in (1) is equivalent to the statement that

$$\langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{y}) - \nabla f(\mathbf{x}) \rangle \geq 0 \quad (2)$$

for all $\mathbf{x}, \mathbf{y} \in \text{dom } f$. This is often called the *monotone gradient* condition.

- (a) Prove that (1) implies (2).
 - (b) (Optional) Prove that (2) implies (1), and thus the conditions are actually equivalent. [Hint: Consider the function $\phi(\alpha) = f(\mathbf{x} + \alpha(\mathbf{y} - \mathbf{x}))$ and recall that the Fundamental Theorem of Calculus tells us that $\phi(1) - \phi(0) = \int_0^1 \phi'(\alpha) d\alpha$.]
 - (c) Consider a one-dimensional differentiable convex function $f(x)$. Assume that $f(x)$ has a unique global minimum x^* . What does the above condition say about $f'(x)$ for $x > x^*$? What about $f'(x)$ for $x < x^*$?
4. **Strong convexity.** In analyzing gradient descent and other algorithms we will sometimes consider the assumption that f is *strongly convex*, i.e., there exists an $m > 0$ such that

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle + \frac{m}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 \quad (3)$$

for all $\mathbf{x}, \mathbf{y} \in \text{dom } f$. Note that above we assume f is differentiable in our definition of strong convexity, and that any f satisfying (3) must be convex. When f is strongly convex, the following statements are equivalent:

- (i) f is strongly convex with constant m , i.e., it satisfies (3).
- (ii) $g(\mathbf{x}) = f(\mathbf{x}) - \frac{m}{2} \|\mathbf{x}\|_2^2$, where $\text{dom } g = \text{dom } f$, is convex.
- (iii) $\langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \rangle \geq m \|\mathbf{x} - \mathbf{y}\|_2^2$ for all $\mathbf{x}, \mathbf{y} \in \text{dom } f$.

In this problem you will show that these are equivalent.

- (a) Prove that (i) is equivalent to (ii).
- (b) Prove that (ii) is equivalent to (iii).

5. **Lipschitz gradients.** In our convergence analysis of gradient descent we will also consider the assumption that ∇f is *Lipschitz*, i.e., there exists an $M > 0$ such that

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq M\|\mathbf{x} - \mathbf{y}\|_2 \quad (4)$$

for all $\mathbf{x}, \mathbf{y} \in \text{dom } f$. When f is convex and differentiable, the following statements are equivalent:

- (i) ∇f is Lipschitz with constant M , i.e., it satisfies (4).
- (ii) $\langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \rangle \leq M\|\mathbf{x} - \mathbf{y}\|_2^2$.
- (iii) $g(\mathbf{x}) = \frac{M}{2}\|\mathbf{x}\|_2^2 - f(\mathbf{x})$ is convex.
- (iv) $f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle + \frac{M}{2}\|\mathbf{y} - \mathbf{x}\|_2^2$.
- (v) $f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle + \frac{1}{2M}\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|_2^2$.
- (vi) $\langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \rangle \geq \frac{1}{M}\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2$.

In the statements above, \mathbf{x}, \mathbf{y} can be arbitrary vectors in $\text{dom } f$. In this problem you will show that these are all equivalent.

- (a) Prove that (i) implies (ii).
- (b) Prove that (ii) implies (iii). [Hint: From problem 3 we know that showing that $g(\mathbf{x})$ is convex is equivalent to establishing the monotonicity of the gradient.]
- (c) Prove that (iii) implies (iv).
- (d) Prove that (iv) implies (v). This is the hardest, so I will give you a short outline of how to approach this. First, use (iv) and what you know about convexity to show that

$$f(\mathbf{x}) - f(\mathbf{y}) \leq \langle \mathbf{x} - \mathbf{z}, \nabla f(\mathbf{x}) \rangle + \langle \mathbf{z} - \mathbf{y}, \nabla f(\mathbf{y}) \rangle + \frac{M}{2}\|\mathbf{z} - \mathbf{y}\|_2^2$$

holds for any $\mathbf{z} \in \text{dom } f$. Since this holds for any \mathbf{z} , you can find the \mathbf{z} that minimizes this upper bound – find this \mathbf{z} by taking the gradient with respect to \mathbf{z} and setting this equal to zero, then plug this back into this expression and simplify.

- (e) Prove that (v) implies (vi).
- (f) Prove that (vi) implies (i).

6. **Convergence Rates.** A central focus when considering different optimization algorithms is the *rate of convergence*. It is often not enough to merely argue that the algorithm converges – we want to know how quickly it will do so, and the rate of convergence allows us to quantify this. In the problems below, we assume that $\{x_k\}$ is a sequence that converges to x^* . We say that $\{x_k\}$ converges *linearly* with a rate of convergence of β if

$$\lim_{k \rightarrow \infty} \frac{|x_{k+1} - x^*|}{|x_k - x^*|} = \beta$$

for some $\beta \in (0, 1)$. If $\beta = 1$ we say that $\{x_k\}$ converges *sublinearly*, and if $\beta = 0$ we say that $\{x_k\}$ converges *superlinearly*.

- (a) Suppose that $\{x_k\}_{k=0}^\infty$ satisfies $|x_{k+1} - x^*| \leq \beta |x_k - x^*|$ for some $0 < \beta < 1$ (and hence converges linearly with rate β). Prove that $|x_k - x^*| \leq \epsilon$ for all

$$k \geq \frac{\log\left(\frac{|x_0 - x^*|}{\epsilon}\right)}{\log \beta}.$$

- (b) Now suppose that $\{x_k\}_{k=0}^\infty$ satisfies $|x_k - x^*| = \frac{k}{k+1} |x_{k-1} - x^*|$. Is this convergence linear, sublinear, or superlinear? How large must k be to ensure that $|x_k - x^*| \leq \epsilon$
- (c) A finer-grained distinction among different kinds of linear/superlinear convergence is the *order of convergence*. We say that $\{x_k\}$ converges with order q if

$$\lim_{k \rightarrow \infty} \frac{|x_{k+1} - x^*|}{|x_k - x^*|^q} = \gamma$$

for some $\gamma > 0$ (not necessarily less than 1). For $q = 1$ this reverts to linear convergence, but $q = 2$ is called *quadratic* convergence, $q = 3$ *cubic* convergence, and so on. Note that q need be an integer. It might not be initially obvious, but quadratic convergence is *much* faster than linear convergence. Consider the two sequences defined by

$$x_k = \frac{1}{2^k} \quad z_k = \frac{1}{2^{2^k}}.$$

Both converge to zero. Show that $\{x_k\}$ converges linearly and compute the rate. Show that $\{z_k\}$ converges quadratically. Submit a plot (on a log scale) that illustrates the difference in how quickly these converge to zero.

7. **Bisection method.** The bisection method is a strategy for one-dimensional problems of the form

$$\underset{x}{\text{minimize}} \quad f(x) \quad \text{subject to} \quad x_l \leq x \leq x_u, \quad (5)$$

where $x_l, x_u \in \mathbb{R}$ satisfy $x_l < x_u$ and $f : \mathbb{R} \rightarrow \mathbb{R}$ is convex. If f is also differentiable, the bisection algorithm can be expressed as follows:

Set initial bounds: $a = x_l, b = x_u$

Initialize $k = 0$

while not converged **do**

$x_k = (a + b)/2$

if $f'(x) > 0$ **then**

$b = x$

else

$a = x$

end if

$k = k + 1$

end while

- (a) Provide an intuitive explanation for why this algorithm makes sense in light of the monotone gradient property of convex functions.
- (b) Assume that f is strictly convex, and hence (5) has a unique solution, denoted x^* . Let x_k denote the estimate provided by the bisection method after k iterations. Argue that x_k converges to x^* .

- (c) Determine whether x_k converges linearly, sublinearly, or superlinearly. If linear, compute the rate of convergence β .
- (d) Propose a variant of the bisection method for solving the unconstrained problem $\text{minimize}_{x \in \mathbb{R}} f(x)$.

8. **Logistic regression.** Logistic regression is a simple, but surprisingly powerful, tool for solving a fundamental problem in machine learning: learning a *classifier* that can distinguish between two classes from a set of training data. Specifically, suppose that we have a set of *training data* $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$, where each $\mathbf{x}_n \in \mathbb{R}^D$ is a “feature vector” and each $y_n \in \{0, 1\}$ is a “label” that indicates which of two classes \mathbf{x}_n corresponds to. The goal in learning a classifier is to find a function h such that $h(\mathbf{x})$ predicts the correct label y for \mathbf{x} that we have never seen before. We can do this by trying to learn a function h that gives us $h(\mathbf{x}_n) = y_n$ on (most of) the training set. In this problem you will explore this idea and implement it on a simple example.

Consider the *logistic* function

$$g(t) = \frac{1}{1 + e^{-t}}.$$

Note that $g(t)$ is always a number between 0 and 1. In logistic regression, we set $t = \mathbf{a}^T \mathbf{x} + b$, where $\mathbf{a} \in \mathbb{R}^D$ and $b \in \mathbb{R}$, so that for any \mathbf{x} we can compute $g(\mathbf{a}^T \mathbf{x} + b)$, returning a number between 0 and 1. Our goal is to learn the parameters \mathbf{a} and b so that we can interpret $g(\mathbf{a}^T \mathbf{x} + b)$ as an estimate of the probability that \mathbf{x} belongs to class 1 (and hence, $1 - g(\mathbf{a}^T \mathbf{x} + b)$ is the probability that \mathbf{x} belongs to class 0). To form a classifier, we then simply compare $g(\mathbf{a}^T \mathbf{x} + b)$ to a threshold in order to estimate the label based on which one has a higher predicted probability, i.e.,

$$h(\mathbf{x}) = \begin{cases} 1 & \text{if } g(\mathbf{a}^T \mathbf{x} + b) \geq \frac{1}{2}, \\ 0 & \text{if } g(\mathbf{a}^T \mathbf{x} + b) < \frac{1}{2}. \end{cases}$$

- (a) Show that $h(\mathbf{x})$ can equivalently be written as

$$h(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{a}^T \mathbf{x} + b \geq 0 \\ 0 & \text{if } \mathbf{a}^T \mathbf{x} + b < 0. \end{cases}$$

- (b) Draw a picture illustrating the set $\{\mathbf{x} \in \mathbb{R}^2 : \mathbf{a}^T \mathbf{x} + b \geq 0\}$ for

$$\mathbf{a} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad b = \frac{1}{2}.$$

- (c) The key question in fitting a logistic regression model is deciding on how to set \mathbf{a} and b based on the training data. To see how this is done, note that if $y_n = 1$, we would like to have $g(\mathbf{a}^T \mathbf{x}_n + b) \approx 1$, in which case

$$\log(g(\mathbf{a}^T \mathbf{x}_n + b)) \approx 0.$$

Similarly, if $y_n = 0$, then we would like $g(\mathbf{a}^T \mathbf{x}_n + b) \approx 0$, or said differently, $1 - g(\mathbf{a}^T \mathbf{x}_n + b) \approx 1$, in which case

$$\log(1 - g(\mathbf{a}^T \mathbf{x}_n + b)) \approx 0.$$

Note that in both cases, the terms on the left-hand side are always negative, so that to make them ≈ 0 corresponds to making these terms as *large* as possible. Combining these desired criteria together, we can express the problem of fitting \mathbf{a} and b as

$$\underset{\mathbf{a}, b}{\text{maximize}} \quad \sum_{n: y_n=1} \log(g(\mathbf{a}^T \mathbf{x}_n + b)) + \sum_{n: y_n=0} \log(1 - g(\mathbf{a}^T \mathbf{x}_n + b)).$$

It is somewhat more convenient to re-express this as the equivalent problem:

$$\underset{\mathbf{a}, b}{\text{maximize}} \quad \sum_{n=1}^N y_n \log(g(\mathbf{a}^T \mathbf{x}_n + b)) + (1 - y_n) \log(1 - g(\mathbf{a}^T \mathbf{x}_n + b)).$$

Show that, by plugging in the formula for g , this problem reduces to

$$\underset{\mathbf{a}, b}{\text{maximize}} \quad \sum_{n=1}^N y_n (\mathbf{a}^T \mathbf{x}_n + b) - \log(1 + e^{\mathbf{a}^T \mathbf{x}_n + b}).$$

- (d) We would now like to use an iterative algorithm like gradient descent to solve this optimization problem. To do this, we will need to calculate the gradient. Note that if we use the notation

$$\tilde{\mathbf{x}}_n = \begin{bmatrix} \mathbf{x}_n \\ 1 \end{bmatrix} \quad \boldsymbol{\theta} = \begin{bmatrix} \mathbf{a} \\ b \end{bmatrix},$$

then we can re-write our objective function as

$$f(\boldsymbol{\theta}) = \sum_{n=1}^N y_n \boldsymbol{\theta}^T \tilde{\mathbf{x}}_n - \log(1 + e^{\boldsymbol{\theta}^T \tilde{\mathbf{x}}_n}).$$

With this notation, compute a formula for the gradient $\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta})$ and the Hessian matrix $\mathbf{D}_f^2(\boldsymbol{\theta})$. [Hint: You may use the answers from homework 1 to avoid any actual computation here, but be careful as the notation in this problem is a little different...]

- (e) Do you think $f(\boldsymbol{\theta})$ is M -smooth for some finite value of M ? Do you think $f(\boldsymbol{\theta})$ is strongly convex? (You do not have to calculate M or m here.)
- (f) Implement a solver for estimating $\boldsymbol{\theta}$ using gradient descent. Test your implementation on the dataset produced by the following code:

```
import numpy as np
from sklearn import datasets
```

```
np.random.seed(2020) # Set random seed so results are repeatable
x, y = datasets.make_blobs(n_samples=100, n_features=2, centers=2, cluster_std=6.0)
```

Use an initial guess of $\boldsymbol{\theta}^{(0)} = \mathbf{0}$. To begin, consider a fixed step size of $\alpha \approx 0.001$ (although you should feel free to play around with α a bit). Report how many iterations your algorithm takes to converge. Remember that you are trying to *minimize* $-f(\boldsymbol{\theta})$. If your algorithm is not converging, a sign error is a likely culprit.

- (g) Now implement the bisection algorithm from the previous problem, and create a modified version of gradient descent in which you update α at each iteration using a line search conducted by running the bisection algorithm to find the optimal α (up to some tolerance that you should select). Report how many iterations are now required when using this strategy to find the “optimal” step size at each iteration. Also report the total number of iterations taken by the combined bisection searches. Does the bisection strategy seem worth the additional costs?