

ELE 538: Large-Scale Optimization Homework 1

Zachary Hervieux-Moore

Wednesday 28th February, 2018

Exercise 1: Strong convexity: Suppose that f is differentiable. Show that the following two statements are equivalent.

i)

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

ii)

$$\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq \mu \|\mathbf{x} - \mathbf{y}\|_2^2, \quad \forall \mathbf{x}, \mathbf{y}$$

Answer:

(i) \Rightarrow (ii):

$$\begin{aligned} f(\mathbf{y}) &\geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}\|_2^2, \quad \forall \mathbf{x}, \mathbf{y} \\ f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle &\geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 \end{aligned}$$

Now adding the norm to both sides,

$$\begin{aligned} &f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 \\ &\geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \mu \|\mathbf{x} - \mathbf{y}\|_2^2 \end{aligned}$$

Now we apply (i) to the LHS of the inequality to get

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

Rearranging yields

$$\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq \mu \|\mathbf{x} - \mathbf{y}\|_2^2, \quad \forall \mathbf{x}, \mathbf{y}$$

(ii) \Rightarrow (i):

We denote the new function $g(\mathbf{x}) = f(\mathbf{x}) - \frac{\mu}{2}\|\mathbf{x}\|_2^2$ and notice that

$$\begin{aligned}\langle \nabla g(\mathbf{x}) - \nabla g(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle &= \langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \mu \langle \mathbf{y} - \mathbf{x}, \mathbf{x} - \mathbf{y} \rangle \\ &= \langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle - \mu \|\mathbf{x} - \mathbf{y}\|_2^2\end{aligned}$$

By (ii) we then have that

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

Therefore, $g(\mathbf{x})$ is convex. Using another notion of convexity,

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

Which leads to

$$\begin{aligned}f(\mathbf{y}) - \frac{\mu}{2}\|\mathbf{y}\|_2^2 &\geq f(\mathbf{x}) - \frac{\mu}{2}\|\mathbf{x}\|_2^2 + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle - \mu \langle \mathbf{x}, \mathbf{y} - \mathbf{x} \rangle \\ f(\mathbf{y}) &\geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\mu}{2}(\|\mathbf{x}\|_2^2 - 2\langle \mathbf{x}, \mathbf{y} \rangle + \|\mathbf{y}\|_2^2) \\ f(\mathbf{y}) &\geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\mu}{2}\|\mathbf{x} - \mathbf{y}\|_2^2, \quad \forall \mathbf{x}, \mathbf{y}\end{aligned}$$

Exercise 2: Subgradients: For each of the following convex functions, explain how to calculate a subgradient at a given $\mathbf{x} = (x_1, \dots, x_n)$.

- a) $f(\mathbf{x}) = \max_{i=1, \dots, m} |\mathbf{a}_i^T \mathbf{x} + b_i|$
- b) $f(\mathbf{x}) = \sup_{0 \leq t \leq 1} p(t)$, where $p(t) = x_1 + x_2 t + \dots + x_n t^{n-1}$
- c) $f(\mathbf{x}) = x_{[1]} + \dots + x_{[k]}$, where $x_{[i]}$ denotes the i^{th} largest element of the vector \mathbf{x}
- d) $f(\mathbf{x}) = \sup_{\mathbf{A}\mathbf{y} \preceq \mathbf{b}} \mathbf{y}^T \mathbf{x}$. (You can assume that the polyhedron defined by $\mathbf{A}\mathbf{y} \preceq \mathbf{b}$ is bounded, where “ \preceq ” denotes component-wise inequality).

Answer:

- a) A subgradient is any $\mathbf{g} = \text{sign}(\mathbf{a}_i^T \mathbf{x} + b_i) \mathbf{a}_i$ such that $f(\mathbf{x}) = |\mathbf{a}_i^T \mathbf{x} + b_i|$ for some $i \in \{1, \dots, m\}$. Suppose \mathbf{a}_k and b_k achieves the max for $f(\mathbf{x})$. Then we check whether or not \mathbf{g} a subgradient. We need to check the following condition:

$$\begin{aligned} f(\mathbf{z}) &\geq f(\mathbf{x}) + \mathbf{g}^T(\mathbf{z} - \mathbf{x}) \\ |\mathbf{a}_k^T \mathbf{z} + b_k| &\geq |\mathbf{a}_k^T \mathbf{x} + b_k| + \text{sign}(\mathbf{a}_k^T \mathbf{x} + b_k) \mathbf{a}_k^T(\mathbf{z} - \mathbf{x}) \end{aligned}$$

From here, it is trivial to check the 4 different cases on the signs of the absolute value terms to verify that this is always true.

- b) First we rewrite the problem as

$$f(x) = \sup_{0 \leq t \leq 1} \mathbf{x}^T \mathbf{a}_t$$

where $\mathbf{a}_t^T = [1 \ t \ \dots \ t^{n-1}]$. By picking the subgradient to be $\mathbf{g} = \mathbf{a}_s$ where s is the argument that maximizes the supremum, we have

$$\begin{aligned} f(\mathbf{z}) &\geq f(\mathbf{x}) + \mathbf{g}^T(\mathbf{z} - \mathbf{x}) \\ f(\mathbf{z}) &\geq f(\mathbf{x}) + \mathbf{a}_s^T \mathbf{z} - \mathbf{a}_s^T \mathbf{x} \\ f(\mathbf{z}) &\geq \mathbf{a}_s^T \mathbf{z} \end{aligned}$$

Where the last line is true because s does not necessarily achieve the supremum at point \mathbf{z} . Thus picking $\mathbf{g} = \mathbf{a}_s$ is infact a subgradient.

So we simply need to find \mathbf{a}_s . By writing down the derivative of the polynomial $p'(t)$ as the characteristic polynomial of a companion matrix, we can efficiently find its eigenvalues using QR decomposition methods. This gives us the critical points of $p(t)$ which we check along with the boundary points $t = 0, 1$. As it is a polynomial on a compact set, one of these points will attain the maximum.

- c) Pick the subgradient to be $\mathbf{g} = \sum_{i=1}^k e_{[i]}$. Where $e_{[i]}$ is the canonical basis vector corresponding to the i^{th} largest entry. Then we have

$$\begin{aligned}
 f(\mathbf{z}) &\geq f(\mathbf{x}) + \mathbf{g}^T(\mathbf{z} - \mathbf{x}) \\
 \sum_{i=1}^k z_{[i]} &\geq \sum_{i=1}^k x_{[i]} + \sum_{i=1}^k e_{[i]}^T(\mathbf{z} - \mathbf{x}) \\
 \sum_{i=1}^k z_{[i]} &\geq \sum_{i=1}^k x_{[i]} + \sum_{i=1}^k (z_{[i]} - x_{[i]}) \\
 &0 \geq 0
 \end{aligned}$$

Thus we do indeed have $\mathbf{g} = \partial f$.

- d) To find a subgradient for $f(\mathbf{x})$, first solve the LP $\sup_{\mathbf{A}\mathbf{y} \preceq \mathbf{b}} \mathbf{y}^T \mathbf{x}$, since it is bounded, the max is achieved. Denote the solutions by \mathbf{y}^* and pick the subgradient to be $\mathbf{g} = \mathbf{y}^*$. Now we verify that this is a gradient.

$$\begin{aligned}
 f(\mathbf{z}) &\geq f(\mathbf{x}) + \mathbf{g}^T(\mathbf{z} - \mathbf{x}) \\
 \sup_{\mathbf{A}\mathbf{w} \preceq \mathbf{b}} \mathbf{w}^T \mathbf{z} &\geq \sup_{\mathbf{A}\mathbf{y} \preceq \mathbf{b}} \mathbf{y}^T \mathbf{x} + \mathbf{y}^{*T}(\mathbf{z} - \mathbf{x}) \\
 \sup_{\mathbf{A}\mathbf{w} \preceq \mathbf{b}} \mathbf{w}^T \mathbf{z} &\geq \mathbf{y}^{*T} \mathbf{x} + \mathbf{y}^{*T} \mathbf{z} - \mathbf{y}^{*T} \mathbf{x} \\
 \sup_{\mathbf{A}\mathbf{w} \preceq \mathbf{b}} \mathbf{w}^T \mathbf{z} &\geq \mathbf{y}^{*T} \mathbf{z}
 \end{aligned}$$

Which is true because \mathbf{y}^* satisfies the constraints and so is feasible to the problem on the LHS.

Exercise 3: A convex function that is not subdifferentiable: Verify that the following function $f : \mathbb{R}_+ \rightarrow \mathbb{R}$ is convex, but not subdifferentiable at $x = 0$:

$$f(x) = \begin{cases} 1, & x = 0 \\ 0, & x > 0 \end{cases}$$

with $\text{domain}(f) = \mathbb{R}_+$

Answer:

First, let us prove f is convex. If $x, y \in (0, \infty)$ and $\lambda \in [0, 1]$, then we have that $f(\lambda x + (1 - \lambda)y) = 0$ and $\lambda f(x) + (1 - \lambda)f(y) = 0$, so it is trivially convex on $(0, \infty)$. Now suppose $x = 0$ and $y \in (0, \infty)$, then we have

$$f(\lambda x + (1 - \lambda)y) = f((1 - \lambda)y) = 0 \leq \lambda f(x) + (1 - \lambda)f(y) = \lambda$$

Thus, $f(x)$ is convex. We will prove that it is not subdifferentiable at $x = 0$ by contradiction. Suppose that there was a subgradient $g \in \mathbb{R}$. Then we would have

$$\begin{aligned} f(z) &\geq f(x) + g(z - x) \quad \forall z \\ f(z) &\geq 1 + gz \quad \forall z \end{aligned}$$

But $f(z) = 0$ for all $z \in (0, \infty)$ so

$$\begin{aligned} 0 &\geq 1 + gz \quad \forall z \\ -gz &\geq 1 \quad \forall z \end{aligned}$$

Note that since $z > 0$, then we must have $g < 0$. Thus picking $z = \frac{1}{-2g}$ which belongs to the interval $(0, \infty)$, leads to a contradiction.

Exercise 4: Matrix norm approximation: We consider the problem of approximating a given matrix $\mathbf{B} \in \mathbb{R}^{p \times q}$ as a linear combination of some other given matrices $\mathbf{A}_i \in \mathbb{R}^{p \times q}$, $i = 1, \dots, n$ as measured by the matrix norm (maximum singular value):

$$\min \|x_1 \mathbf{A}_1 + \dots + x_n \mathbf{A}_n - \mathbf{B}\|_2$$

- a) Explain how to find a subgradient of the objective function at \mathbf{x}
- b) Generate a random instance of the problem with $n = 5, p = 3, q = 6$. Use CVX to find the optimal value of f^* of the problem. Use a subgradient method to solve the problem, starting from $\mathbf{x} = \mathbf{0}$. Plot $f - f^*$ versus iteration. Experiment with several step size sequences.

Answer:

- a) We note the the matrix norm of a matrix C is equivalent to solving the following problem

$$\sup_{\|y\|_2 \leq 1} y^T C^T C y$$

Or put more succinctly, the square root of the largest eigenvalue of $C^T C$. In our case, $C = x_1 \mathbf{A}_1 + \dots + x_n \mathbf{A}_n - \mathbf{B}$. Thus, $C^T C$ for us is

$$C^T C = \sum_{i=1}^n \sum_{j=1}^n x_i A_i^T A_j x_j - \sum_{i=1}^n (x_i A_i^T B + x_i B^T A_i) + B^T B$$

Suppose \mathbf{y} is the eigenvector corresponding to the largest singular value, then we have our objective being equal to

$$f(\mathbf{x}) = \sum_{i=1}^n \sum_{j=1}^n x_i \mathbf{y}^T A_i^T A_j \mathbf{y} x_j - \sum_{i=1}^n (x_i \mathbf{y}^T A_i^T B \mathbf{y} + x_i \mathbf{y}^T B^T A_i \mathbf{y}) + \mathbf{y}^T B^T B \mathbf{y}$$

Thus, a subgradient of $f(\mathbf{x})$ can be obtained by taking the gradient of the above

$$\mathbf{g}_i = \mathbf{y}^T \left(2x_i A_i^T A_i + \sum_{i \neq j} (x_j A_i^T A_j + x_j A_j^T A_i) - A_i^T B - B^T A_i \right) \mathbf{y}$$

- b) The code is appended below but the following three figures show the convergence for difference step size schemes.

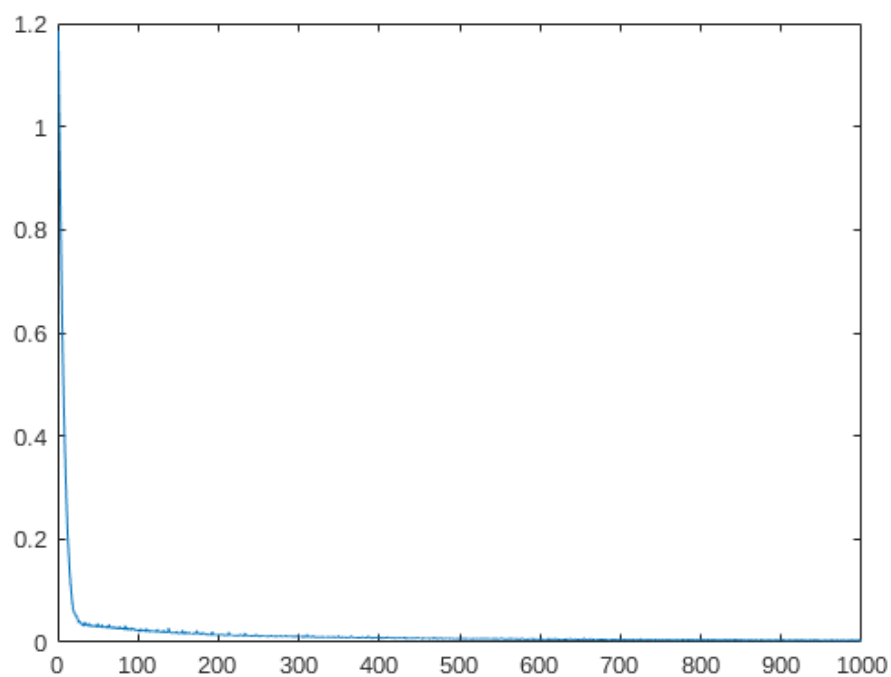


Figure 1: Difference between $f(x_t)$ and f^{opt} at each iteration using a Polyak step size.

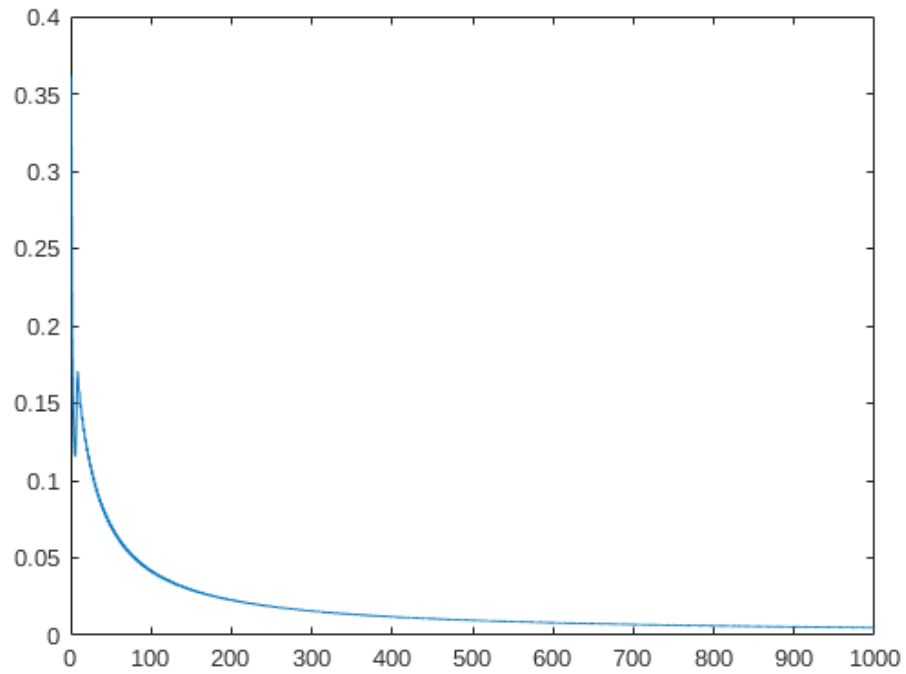


Figure 2: Difference between $f(x_t)$ and f^{opt} at each iteration using a $1/t$ step size. Note the slower convergence.

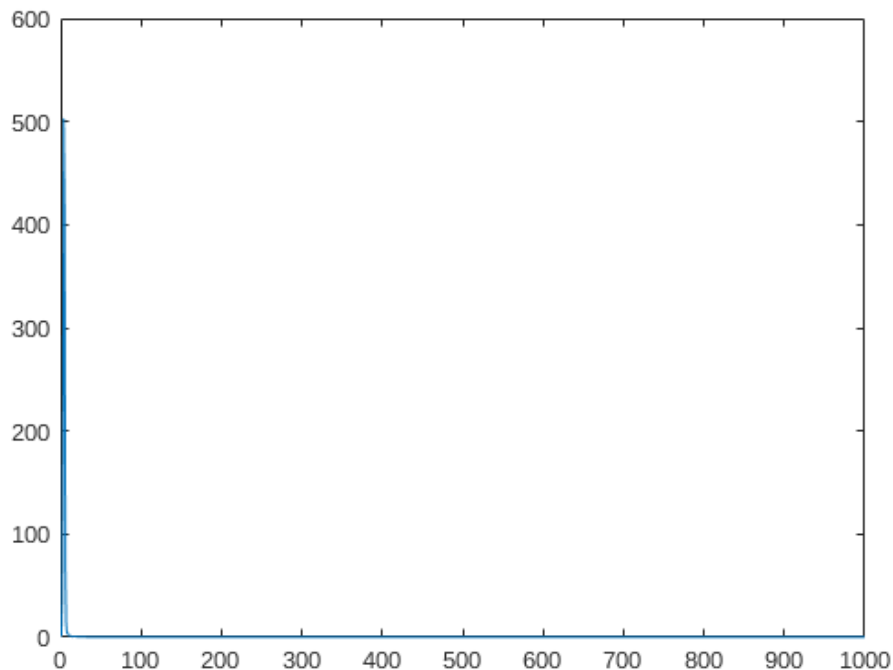


Figure 3: Difference between $f(x_t)$ and f^{opt} at each iteration using a $1/t^2$ step size. Note that divergence at the beginning when the step size is too big. The convergence is also much slower than the other two but the scale of the y-axis is too large to tell.

Code:

```
clear;
clc;

n = 5; p = 3; q = 6;
A = randn(p,q,n);
B = randn(p,q);

cvx_begin
    variable x(n)
    minimize( norm(A(:, :, 1)*x(1)+A(:, :, 2)*x(2)+A(:, :, 3)*x(3)+A
        (:, :, 4)*x(4)+A(:, :, 5)*x(5) - B, 2) )
cvx_end

f_opt = cvx_optval;
```

```

% Subgradient method, no need to project as we are in  $\mathbb{R}^n$ 
T = 1000;
y = zeros(T,5);
g = zeros(T,5);
f = zeros(T,1);
tmp = zeros(6,6);

for t = 1:T
    inner_mat = A(:, :, 1)*y(t,1)+A(:, :, 2)*y(t,2)+A(:, :, 3)*y(t,3)+A(:, :, 4)*y(t,4)+A(:, :, 5)*y(t,5) - B;
    f(t) = norm(inner_mat, 2);
    [U,S,V] = svds(inner_mat, 1);

    % Calculate subgradient
    for i = 1:n
        tmp = 2*y(t,i)*A(:, :, i)'*A(:, :, i);

        for j = 1:n
            if j ~= i
                tmp = tmp + y(t,j)*A(:, :, i)'*A(:, :, j) + y(t,j)*A(:, :, j)'*A(:, :, i);
            end
        end

        tmp = tmp - A(:, :, i)'*B - B'*A(:, :, i);
        g(t,i) = V'*tmp*V;
    end

    eta = (f(t)-f_opt)/(norm(g(t, :), 2)^2);
    %eta = 1/t^2;
    %eta = 1/(t+30);
    y(t+1,:) = y(t, :) - eta*g(t, :);
end

plot(f-f_opt)
f(T)-f_opt

```

Exercise 5: Step sizes that guarantee moving closer to the optimal

set: Consider the subgradient method iteration $\mathbf{x}^+ = \mathbf{x} - \eta \mathbf{g}$, where $\mathbf{g} \in \partial f(\mathbf{x})$. Let f^* be the optimal objective value. Show that if $\eta < \frac{2(f(\mathbf{x}) - f^*)}{\|\mathbf{g}\|_2^2}$ (which is twice Polyak's optimal step size value) we have

$$\|\mathbf{x}^+ - \mathbf{x}^*\|_2 < \|\mathbf{x} - \mathbf{x}^*\|_2$$

for any optimal point \mathbf{x}^* .

Remark: Methods in which successive iterates move closer to the optimal set are called *Fejer monotone*. Thus, the subgradient method, with Polyak's optimal step size, is *Fejer monotone*.

Answer:

The proof follows immediately using the majorizing function presented by Lemma 4.1 in the notes. We have by the lemma:

$$\|\mathbf{x}^+ - \mathbf{x}^*\|_2^2 \leq \|\mathbf{x} - \mathbf{x}^*\|_2^2 - 2\eta(f(\mathbf{x}) - f^*) + \eta^2\|\mathbf{g}\|_2^2$$

If we wish to have the result that $\|\mathbf{x}^+ - \mathbf{x}^*\|_2 < \|\mathbf{x} - \mathbf{x}^*\|_2$, then we must have that

$$-2\eta(f(\mathbf{x}) - f^*) + \eta^2\|\mathbf{g}\|_2^2 < 0$$

Simple rearranging of this yields

$$\eta < \frac{2(f(\mathbf{x}) - f^*)}{\|\mathbf{g}\|_2^2}$$

Note that dividing by η is justified as if $\eta = 0$ then we are at the optimal point.

Exercise 6: Gradient of squared distance (bonus): Define the Euclidean projection onto a closed convex set \mathcal{C} as

$$\mathcal{P}_{\mathcal{C}}(\mathbf{x}) := \arg \min_{\mathbf{z} \in \mathcal{C}} \|\mathbf{z} - \mathbf{x}\|_2$$

and let

$$\text{dist}_{\mathcal{C}}(\mathbf{x}) := \|\mathbf{x} - \mathcal{P}_{\mathcal{C}}(\mathbf{x})\|_2$$

Show that the gradient of the squared distance $f(\mathbf{x}) := \frac{1}{2} \text{dist}_{\mathcal{C}}^2(\mathbf{x})$ is

$$\nabla f(\mathbf{x}) = \mathbf{x} - \mathcal{P}_{\mathcal{C}}(\mathbf{x})$$

Here, you can assume (without proof) that $f(\mathbf{x})$ is convex.

Answer:

For conciseness, we use the shorthand notation $\mathcal{P}_{\mathcal{C}}(\mathbf{x}) = \mathcal{P}_{\mathbf{x}}$. To show that the $\mathbf{x} - \mathcal{P}_{\mathbf{x}}$ is the gradient of $f(\mathbf{x})$ we must show

$$f(\mathbf{y}) - f(\mathbf{x}) \geq (\mathbf{y} - \mathbf{x})^T (\mathbf{x} - \mathcal{P}_{\mathbf{x}})$$

Plugging in the definitions into the above

$$\|\mathbf{y} - \mathcal{P}_{\mathbf{y}}\|_2^2 - \|\mathbf{x} - \mathcal{P}_{\mathbf{x}}\|_2^2 \geq 2(\mathbf{y} - \mathbf{x})^T (\mathbf{x} - \mathcal{P}_{\mathbf{x}})$$

Now we introduce some terms to the RHS and simplify

$$\begin{aligned} \|\mathbf{y} - \mathcal{P}_{\mathbf{y}}\|_2^2 - \|\mathbf{x} - \mathcal{P}_{\mathbf{x}}\|_2^2 &\geq 2(\mathbf{y} - \mathcal{P}_{\mathbf{y}} + \mathcal{P}_{\mathbf{y}} - \mathbf{x})^T (\mathbf{x} - \mathcal{P}_{\mathbf{x}}) \\ \|\mathbf{y} - \mathcal{P}_{\mathbf{y}}\|_2^2 - \|\mathbf{x} - \mathcal{P}_{\mathbf{x}}\|_2^2 &\geq 2(\mathbf{y} - \mathcal{P}_{\mathbf{y}})^T (\mathbf{x} - \mathcal{P}_{\mathbf{x}}) + 2(\mathcal{P}_{\mathbf{y}} - \mathbf{x})^T (\mathbf{x} - \mathcal{P}_{\mathbf{x}}) \\ \|\mathbf{y} - \mathcal{P}_{\mathbf{y}}\|_2^2 + \|\mathbf{x} - \mathcal{P}_{\mathbf{x}}\|_2^2 - 2(\mathbf{y} - \mathcal{P}_{\mathbf{y}})^T (\mathbf{x} - \mathcal{P}_{\mathbf{x}}) &\geq 2(\mathcal{P}_{\mathbf{y}} - \mathbf{x})^T (\mathbf{x} - \mathcal{P}_{\mathbf{x}}) + 2\|\mathbf{x} - \mathcal{P}_{\mathbf{x}}\|_2^2 \\ \|\mathbf{y} - \mathcal{P}_{\mathbf{y}} - \mathbf{x} + \mathcal{P}_{\mathbf{x}}\|_2^2 &\geq 2(\mathcal{P}_{\mathbf{y}} - \mathbf{x})^T (\mathbf{x} - \mathcal{P}_{\mathbf{x}}) + 2\|\mathbf{x} - \mathcal{P}_{\mathbf{x}}\|_2^2 \end{aligned}$$

Now, we expand the norm on the LHS differently

$$\|\mathbf{y} - \mathbf{x}\|_2^2 + 2(\mathbf{y} - \mathbf{x})^T(\mathcal{P}_x - \mathcal{P}_y) + \|\mathcal{P}_x - \mathcal{P}_y\|_2^2 \geq 2(\mathcal{P}_y - \mathbf{x})^T(\mathbf{x} - \mathcal{P}_x) + 2\|\mathbf{x} - \mathcal{P}_x\|_2^2$$

Using the non-expansiveness property of projections, we can upper bound the LHS

$$2\|\mathbf{y} - \mathbf{x}\|_2^2 + 2(\mathbf{y} - \mathbf{x})^T(\mathcal{P}_x - \mathcal{P}_y) \geq 2(\mathcal{P}_y - \mathbf{x})^T(\mathbf{x} - \mathcal{P}_x) + 2\|\mathbf{x} - \mathcal{P}_x\|_2^2$$

Simplifying this expression

$$\begin{aligned} \|\mathbf{y} - \mathbf{x}\|_2^2 + (\mathbf{y} - \mathbf{x})^T(\mathcal{P}_x - \mathcal{P}_y) &\geq (\mathcal{P}_y - \mathcal{P}_x)^T(\mathbf{x} - \mathcal{P}_x) \\ \|\mathbf{y} - \mathbf{x}\|_2^2 &\geq (\mathcal{P}_y - \mathcal{P}_x)^T(\mathbf{x} - \mathcal{P}_x + \mathbf{y} - \mathbf{x}) \\ \|\mathbf{y} - \mathbf{x}\|_2^2 &\geq (\mathcal{P}_y - \mathcal{P}_x)^T(\mathbf{y} - \mathcal{P}_x) \end{aligned}$$

Introducing $\mathcal{P}_y - \mathcal{P}_y$ to the inner product of the RHS yields

$$\|\mathbf{y} - \mathbf{x}\|_2^2 \geq \|\mathcal{P}_y - \mathcal{P}_x\|_2^2 + (\mathcal{P}_y - \mathcal{P}_x)^T(\mathbf{y} - \mathcal{P}_y)$$

Now, we note that the second term in the RHS $(\mathcal{P}_y - \mathcal{P}_x)^T(\mathbf{y} - \mathcal{P}_y) \geq 0$ by the convexity of \mathcal{C} . To see this, note that $\mathbf{y} - \mathcal{P}_y$ is perpendicular to the tangent at \mathcal{P}_y and that \mathcal{P}_x and \mathcal{P}_y are contained inside \mathcal{C} and so these two vectors must have a positive inner product. Thus, we can lower bound the LHS by dropping this term

$$\|\mathbf{y} - \mathbf{x}\|_2^2 \geq \|\mathcal{P}_y - \mathcal{P}_x\|_2^2$$

Which is our non-expansiveness property that we know is true. Thus, our original inequality must be true and we conclude that

$$\nabla f(\mathbf{x}) = \mathbf{x} - \mathcal{P}_{\mathcal{C}}(\mathbf{x})$$