

1 Lagrange Multipliers and Convex Duality Recap

Recall the convex optimization problem we studied last lecture,

$$\begin{aligned} \min \quad & \mathcal{E}(\mathbf{y}) \\ \text{s.t.} \quad & \mathbf{A}\mathbf{y} = \mathbf{b} \\ & \mathbf{c}(\mathbf{y}) \leq \mathbf{0}, \end{aligned} \tag{1}$$

where $\mathcal{E}(\mathbf{y}) : S \rightarrow \mathbb{R}$ is defined on a subset $S \subseteq \mathbb{R}^n$, $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{c}(\mathbf{y})$ is a vector of constraints $\mathbf{c}(\mathbf{y}) = (c_i(\mathbf{y}))_{i \in [k]}$. For every $i \in [k]$ the function $c_i : S \rightarrow \mathbb{R}$ is convex. We call (1) the primal (problem) and denote its optimal value by α^* .

The associated Lagrangian is defined by

$$L(\mathbf{y}, \mathbf{x}, \mathbf{s}) = \mathcal{E}(\mathbf{y}) + \mathbf{x}^T(\mathbf{b} - \mathbf{A}\mathbf{y}) + \mathbf{s}^T \mathbf{c}(\mathbf{y}).$$

where $\mathbf{x} \in \mathbb{R}^m$, $\mathbf{s} \in \mathbb{R}^k$ are dual variables. The dual (problem) is given by

$$\max_{\substack{\mathbf{x}, \mathbf{s} \\ \mathbf{s} \geq \mathbf{0}}} L(\mathbf{x}, \mathbf{s}) \tag{2}$$

whose optimal value is denoted by β^* . The dual is always a convex optimization problem even though the primal is non-convex. The optimal value of the primal (1) can also be written as

$$\alpha^* = \inf_{\mathbf{y}} \sup_{\mathbf{x}; \mathbf{s} \geq \mathbf{0}} L(\mathbf{y}, \mathbf{x}, \mathbf{s}), \tag{3}$$

where no constraint is imposed on the primal variable \mathbf{y} . The optimal value of the dual (2) is

$$\beta^* = \sup_{\mathbf{x}; \mathbf{s} \geq \mathbf{0}} \inf_{\mathbf{y}} L(\mathbf{y}, \mathbf{x}, \mathbf{s}). \tag{4}$$

Note the only difference between (3) and (4) is that the positions of “inf” and “sup” are swapped. The weak duality theorem states that the dual optimal value is a lower bound of the primal optimal value, i.e. $\beta^* \leq \alpha^*$.

The Slater's condition for (1) requires the existence of a *strictly feasible* point, i.e. there exists $\tilde{\mathbf{y}} \in S$ s.t. $\mathbf{A}\tilde{\mathbf{y}} = \mathbf{b}$ and $\mathbf{c}(\tilde{\mathbf{y}}) < \mathbf{0}$. This means that the strictly feasible point $\tilde{\mathbf{y}}$ lies inside the interior of the set $\{\mathbf{y} : \mathbf{c}(\mathbf{y}) \leq \mathbf{0}\}$ defined by the inequality constraints. The strong duality theorem says, the Slater's condition implies strong duality, $\beta^* = \alpha^*$.

Example. In lecture 10, we gave a combinatorial proof of the min-cut max-flow theorem, and showed that the min-cut problem can be expressed as a linear program. Now, we will use the

strong duality theorem to give an alternative proof, and directly find the min-cut linear program is the dual program to our maximum flow linear program.

We will assume that Slater's condition holds for our primal program. Since scaling the flow down enough will always ensure that capacity constraints are strictly satisfied i.e. $\mathbf{f} < \mathbf{c}$, the only concern is to make sure that non-negativity constraints are satisfied. This means that there is an s - t flow that sends a non-zero flow on every edge. In fact, this may not always be possible, but it is easy to detect such edges and remove them without changing the value of the program: an edge (u, v) should be removed if there is no path s to u or no path v to t . We can identify all such edges using a BFS from s along the directed edges and a BFS along reversed directed edges from t .

Slater's condition holds whenever there is a directed path from s to t with non-zero capacity (and if there is not, the maximum flow and minimum cut are both zero).

$$\begin{aligned}
& \min_{\substack{F \in \mathbb{R} \\ \mathbf{B}\mathbf{f} = F\mathbf{b}_{s,t} \\ \mathbf{0} \leq \mathbf{f} \leq \mathbf{c}}} -F = \min_{F; \mathbf{f} \geq \mathbf{0}} \max_{\mathbf{x}; \mathbf{s} \geq \mathbf{0}} -F + \mathbf{x}^\top (F\mathbf{b}_{s,t} - \mathbf{B}\mathbf{f}) + (\mathbf{f} - \mathbf{c})^\top \mathbf{s} \\
& \quad \text{(Slater's condition } \implies \text{ strong duality)} \\
- & \max_{\substack{F \in \mathbb{R} \\ \mathbf{B}\mathbf{f} = F\mathbf{b}_{s,t} \\ \mathbf{0} \leq \mathbf{f} \leq \mathbf{c}}} F = \max_{\mathbf{x}; \mathbf{s} \geq \mathbf{0}} \min_{F; \mathbf{f} \geq \mathbf{0}} F(\mathbf{b}_{s,t}^\top \mathbf{x} - 1) + \mathbf{f}^\top (\mathbf{s} - \mathbf{B}^\top \mathbf{x}) - \mathbf{c}^\top \mathbf{s} \\
& = \max_{\substack{\mathbf{x}; \mathbf{s} \geq \mathbf{0} \\ \mathbf{b}_{s,t}^\top \mathbf{x} = 1 \\ \mathbf{s} \geq \mathbf{B}^\top \mathbf{x}}} -\mathbf{c}^\top \mathbf{s} \\
& \max_{\substack{F \in \mathbb{R} \\ \mathbf{B}\mathbf{f} = F\mathbf{b}_{s,t} \\ \mathbf{0} \leq \mathbf{f} \leq \mathbf{c}}} F = \min_{\substack{\mathbf{x}; \mathbf{s} \geq \mathbf{0} \\ \mathbf{b}_{s,t}^\top \mathbf{x} = 1 \\ \mathbf{s} \geq \mathbf{B}^\top \mathbf{x}}} \mathbf{c}^\top \mathbf{s} \tag{5}
\end{aligned}$$

The LHS of (5) is exactly the LP formulation of max-flow, while the RHS is exactly the LP formulation of min-cut.

Let \mathbf{y}^* be an optimizer of the primal problem and $\mathbf{x}^*, \mathbf{s}^*$ for the dual. Since $\mathbf{y}^*, \mathbf{x}^*, \mathbf{s}^*$ are also primal/dual feasible, then

$$\alpha^* = \mathcal{E}(\mathbf{y}^*) \geq L(\mathbf{y}^*, \mathbf{x}^*, \mathbf{s}^*) \geq L(\mathbf{x}^*, \mathbf{s}^*) = \beta^*.$$

Now, suppose strong duality holds, i.e. $\alpha^* = \beta^*$. Then,

$$\mathcal{E}(\mathbf{y}^*) = L(\mathbf{y}^*, \mathbf{x}^*, \mathbf{s}^*) = L(\mathbf{x}^*, \mathbf{s}^*).$$

If further assume $\mathcal{E} : S \rightarrow \mathbb{R}$ and \mathbf{c} are differentiable and the minimizer \mathbf{y}^* is not on the boundary of S , then

$$\begin{aligned}
& \nabla_{\mathbf{y}} L(\mathbf{y}, \mathbf{x}^*, \mathbf{s}^*) \Big|_{\mathbf{y}=\mathbf{y}^*} = \mathbf{0} \\
& \nabla_{\mathbf{y}} \left(\mathcal{E}(\mathbf{y}) + (\mathbf{x}^*)^\top (\mathbf{b} - \mathbf{A}\mathbf{y}) + (\mathbf{s}^*)^\top \mathbf{c}(\mathbf{y}) \right) \Big|_{\mathbf{y}=\mathbf{y}^*} = \mathbf{0} \\
& \nabla_{\mathbf{y}} \mathcal{E}(\mathbf{y}^*) - \mathbf{A}^\top \mathbf{x}^* - \mathbf{s}^{\top} \nabla \mathbf{c}(\mathbf{y}^*) = \mathbf{0}.
\end{aligned}$$

And this connects to our point of the parallel gradients.

Furthermore, we also have *complementary slackness*,

$$\mathbf{s}^*(i) \cdot \mathbf{c}_i(\mathbf{y}^*) = 0 \text{ for all } i,$$

which means the i -th optimal dual variable $\mathbf{s}^*(i)$ is zero unless the i -th constraint is active at the optimum, i.e. $\mathbf{c}_i(\mathbf{y}^*) = 0$.

2 Karush-Kuhn-Tucker Theorem

Let us continue to consider the convex optimization problem (1) with \mathcal{E} and \mathbf{c} differentiable and assume that the domain S of \mathcal{E} is an *open* convex set. We assume S is open to rule out having an optimal \mathbf{y} on the boundary of S . Suppose $\tilde{\mathbf{x}}, \tilde{\mathbf{y}}, \tilde{\mathbf{s}}$ satisfy the following conditions:

- $\mathbf{A}\tilde{\mathbf{y}} = \mathbf{b}$ and $\mathbf{c}(\tilde{\mathbf{y}}) \leq 0$ (primal feasible)
- $\tilde{\mathbf{s}} \geq 0$ (dual feasible)
- $\nabla_{\mathbf{y}}\mathcal{E}(\tilde{\mathbf{y}}) - \mathbf{A}^\top \tilde{\mathbf{x}} - \tilde{\mathbf{s}}^\top \nabla \mathbf{c}(\mathbf{y}^*) = \mathbf{0}$ ($\nabla_{\mathbf{y}}L(\tilde{\mathbf{y}}, \tilde{\mathbf{x}}, \tilde{\mathbf{s}}) = \mathbf{0}$)
- $\tilde{\mathbf{s}}(i) \cdot \mathbf{c}_i(\tilde{\mathbf{y}}) = 0$ for all i (complementary slackness)

These conditions are called the *Karush-Kuhn-Tucker* (KKT) conditions.

Theorem 2.1. $\tilde{\mathbf{x}}, \tilde{\mathbf{y}}, \tilde{\mathbf{s}}$ are primal/dual optimal.

Proof. $\tilde{\mathbf{y}}$ is global minimizer of $\mathbf{y} \mapsto L(\mathbf{y}, \tilde{\mathbf{x}}, \tilde{\mathbf{s}})$, since this function is convex with vanishing gradient at $\tilde{\mathbf{y}}$. Hence,

$$L(\tilde{\mathbf{y}}, \tilde{\mathbf{x}}, \tilde{\mathbf{s}}) = \inf_{\mathbf{y}} L(\mathbf{y}, \tilde{\mathbf{x}}, \tilde{\mathbf{s}}) = L(\tilde{\mathbf{x}}, \tilde{\mathbf{s}}) \leq \beta^*.$$

On the other hand, due to primal feasibility and complementary slackness,

$$L(\tilde{\mathbf{y}}, \tilde{\mathbf{x}}, \tilde{\mathbf{s}}) = \mathcal{E}(\tilde{\mathbf{y}}) + \tilde{\mathbf{x}}^\top (\mathbf{b} - \mathbf{A}^\top \tilde{\mathbf{y}}) + \tilde{\mathbf{s}}^\top \mathbf{c}(\tilde{\mathbf{y}}) = \mathcal{E}(\tilde{\mathbf{y}}) \geq \alpha^*.$$

Thus, $\beta^* \geq \alpha^*$. But also $\beta^* \leq \alpha^*$ by weak duality. Therefore, $\beta^* = \alpha^*$ and $\tilde{\mathbf{y}}, \tilde{\mathbf{x}}, \tilde{\mathbf{s}}$ are primal/dual optimal. \square

3 Fenchel Conjugate

Definition 3.1 (Fenchel conjugate). Given a (convex) function $\mathcal{E} : S \subseteq \mathbb{R}^n \rightarrow \mathbb{R}$, its *Fenchel conjugate* is a function $\mathcal{E}^* : \mathbb{R}^n \rightarrow \mathbb{R}$ defined as

$$\mathcal{E}^*(\mathbf{z}) = \sup_{\mathbf{y} \in S} \langle \mathbf{z}, \mathbf{y} \rangle - \mathcal{E}(\mathbf{y}).$$

Remark 3.2. \mathcal{E}^* is a convex function whether \mathcal{E} is convex or not, since $\mathcal{E}^*(z)$ is pointwise supremum of a family of convex (here, affine) functions of z .

In this course, we have only considered convex function that are real-valued and continuous and defined on a convex domain. For any such \mathcal{E} , we have $\mathcal{E}^{**} = \mathcal{E}$, i.e. the Fenchel conjugate of the Fenchel conjugate is the original function. This is a consequence of the Fenchel-Moreau theorem, which establishes this under slightly more general conditions. We will not prove this generally, but as part of Theorem 3.3 below, we prove it under more restrictive assumptions.

Example. Let $\mathcal{E}(\mathbf{y}) = \frac{1}{p} \|\mathbf{y}\|_p^p$ ($p > 1$). We want to evaluate its Fenchel conjugate \mathcal{E}^* at any given point $\mathbf{z} \in \mathbb{R}^n$. Since \mathcal{E} is convex and differentiable, the supremum must be achieved at some \mathbf{y} with vanishing gradient

$$\nabla_{\mathbf{y}} \langle \mathbf{z}, \mathbf{y}^* \rangle - \nabla \mathcal{E}(\mathbf{y}^*) = \mathbf{z} - \nabla \mathcal{E}(\mathbf{y}^*) = \mathbf{0} \iff \mathbf{z} = \nabla \mathcal{E}(\mathbf{y}^*).$$

It's not difficult to see, for all i ,

$$z(i) = \text{sgn}(\mathbf{y}(i)) |\mathbf{y}(i)|^{p-1}.$$

Then,

$$\begin{aligned} \mathcal{E}^*(\mathbf{z}) &= \langle \mathbf{z}, \mathbf{y} \rangle - \mathcal{E}(\mathbf{y}) \\ &= \sum_i |z(i)|^{\frac{1}{p-1}+1} - \frac{1}{p} |z(i)|^{\frac{p}{p-1}} \\ &\text{(define } q \text{ s.t. } \frac{1}{q} + \frac{1}{p} = 1) \\ &= \frac{1}{q} \|\mathbf{z}\|_q^q \end{aligned}$$

More generally, given a convex and differentiable function $\mathcal{E} : S \rightarrow \mathbb{R}$, if there exists $\mathbf{y} \in S$ s.t. $\mathbf{z} = \nabla \mathcal{E}(\mathbf{y})$, then $\mathcal{E}^*(\mathbf{z}) = (\mathbf{y}^*)^\top \nabla \mathcal{E}(\mathbf{y}^*) - \mathcal{E}(\mathbf{y}^*)$. The Fenchel conjugate and Lagrange duality are closely related, which is demonstrated in the following example.

Example. Consider a convex optimization problem with only linear constraints,

$$\begin{aligned} \min_{\mathbf{y} \in \mathbb{R}^n} \quad & \mathcal{E}(\mathbf{y}) \\ \text{s.t.} \quad & \mathbf{A}\mathbf{y} = \mathbf{b} \end{aligned}$$

where $\mathcal{E} : \mathbb{R}^n \rightarrow \mathbb{R}$ is a convex function and $\mathbf{A} \in \mathbb{R}^{m \times n}$. Then the corresponding dual problem is

$$\begin{aligned} \sup_{\mathbf{x} \in \mathbb{R}^m} \inf_{\mathbf{y} \in \mathbb{R}^n} \mathcal{E}(\mathbf{y}) + \mathbf{x}^\top (\mathbf{b} - \mathbf{A}\mathbf{y}) &= \sup_{\mathbf{x} \in \mathbb{R}^m} \mathbf{b}^\top \mathbf{x} - \sup_{\mathbf{y} \in \mathbb{R}^n} (\mathbf{x}^\top \mathbf{A}\mathbf{y} - \mathcal{E}(\mathbf{y})) \\ &= \sup_{\mathbf{x} \in \mathbb{R}^m} \mathbf{b}^\top \mathbf{x} - \mathcal{E}^*(\mathbf{A}^\top \mathbf{x}) \end{aligned}$$

Theorem 3.3 (properties of Fenchel conjugate). *When $\mathcal{E} : \mathbb{R}^n \rightarrow \mathbb{R}$ is bounded below, strictly convex, differentiable, with a Hessian that is positive definite everywhere, we have the following three properties:*

1. $\nabla \mathcal{E}(\nabla \mathcal{E}^*(\mathbf{z})) = \mathbf{z}$ and $\nabla \mathcal{E}^*(\nabla \mathcal{E}(\mathbf{y})) = \mathbf{y}$
2. $(\mathcal{E}^*)^* = \mathcal{E}$, i.e. the Fenchel conjugate of the Fenchel conjugate is the original function.

$$3. \mathbf{H}_{\mathcal{E}^*}(\nabla \mathcal{E}(\mathbf{y})) = \mathbf{H}_{\mathcal{E}}^{-1}(\mathbf{y})$$

$$\begin{array}{ccc}
\text{primal point } \mathbf{y} & \begin{array}{c} \xrightarrow{\nabla_{\mathbf{y}} \mathcal{E}} \\ \xleftarrow{\nabla_{\mathbf{z}} \mathcal{E}^*} \end{array} & \text{dual point } \mathbf{z} \\
\\
\text{gradient } \nabla_{\mathbf{y}} \mathcal{E}(\mathbf{y}) = \mathbf{z} & & \text{gradient } \nabla_{\mathbf{z}} \mathcal{E}^*(\mathbf{z}) = \mathbf{y} \\
\\
\text{Hessian } \mathbf{H}_{\mathcal{E}}(\mathbf{y}) = \mathbf{H}_{\mathcal{E}^*}^{-1}(\mathbf{z}) & & \text{Hessian } \mathbf{H}_{\mathcal{E}^*}(\mathbf{z}) = \mathbf{H}_{\mathcal{E}}^{-1}(\mathbf{y})
\end{array}$$

Figure 1: Properties of Fenchel conjugate

Proof sketch. Part 1. Given \mathbf{z} , let $\mathbf{y} = \mathbf{y}(\mathbf{z})$ be the \mathbf{y} achieving the supremum in the Fenchel conjugate program such that $\mathcal{E}^*(\mathbf{z}) = \langle \mathbf{z}, \mathbf{y}(\mathbf{z}) \rangle - \mathcal{E}(\mathbf{y}(\mathbf{z}))$. One can show that because \mathcal{E} is bounded below and strictly convex, $\mathbf{y} \mapsto \langle \mathbf{z}, \mathbf{y} \rangle - \mathcal{E}(\mathbf{y})$ is bounded above and is strictly concave, and is hence maximized at some \mathbf{y} such that $\mathbf{z} = \nabla \mathcal{E}(\mathbf{y}(\mathbf{z}))$. Since \mathcal{E} is strictly convex, $\mathbf{y}(\mathbf{z})$ is unique. Then, using the product rule and composition rule of derivatives,

$$\begin{aligned}
\nabla \mathcal{E}^*(\mathbf{z}) &= \nabla_{\mathbf{z}} (\langle \mathbf{z}, \mathbf{y}(\mathbf{z}) \rangle - \mathcal{E}(\mathbf{y}(\mathbf{z}))) \\
&= \mathbf{y}(\mathbf{z}) + \text{diag}(\mathbf{z}) \nabla_{\mathbf{z}} \mathbf{y}(\mathbf{z}) - \text{diag}(\underbrace{\nabla_{\mathbf{y}} \mathcal{E}(\mathbf{y}(\mathbf{z}))}_{=\mathbf{z}}) \nabla_{\mathbf{z}} \mathbf{y}(\mathbf{z}) \\
&= \mathbf{y}(\mathbf{z})
\end{aligned}$$

Thus we have $\nabla_{\mathbf{y}} \mathcal{E}(\mathbf{y}(\mathbf{z})) = \mathbf{z}$ and $\nabla_{\mathbf{z}} \mathcal{E}^*(\mathbf{z}) = \mathbf{y}(\mathbf{z})$. Combining the two, we have $\nabla \mathcal{E}(\nabla \mathcal{E}^*(\mathbf{z})) = \mathbf{z}$.

We can also see that for any \mathbf{y} , there exists a \mathbf{z} such that $\nabla_{\mathbf{z}} \mathcal{E}^*(\mathbf{z}) = \mathbf{y}$, namely, this is attained by $\mathbf{z} = \nabla_{\mathbf{y}} \mathcal{E}(\mathbf{y})$. Thus, $\nabla \mathcal{E}^*(\nabla \mathcal{E}(\mathbf{y})) = \mathbf{y}$.

Part 2. Observe that

$$\mathcal{E}^{**}(\mathbf{u}) = \sup_{\mathbf{z} \in \mathbb{R}^n} \langle \mathbf{u}, \mathbf{z} \rangle - \mathcal{E}^*(\mathbf{z})$$

and let $\mathbf{z}(\mathbf{u})$ denote the \mathbf{z} obtaining the supremum, in the above program. We then have $\mathbf{u} = \nabla \mathcal{E}^*(\mathbf{z}(\mathbf{u}))$. Letting $\mathbf{y}(\mathbf{z})$ be defined as in Part 1, we get $\mathbf{y}(\mathbf{z}(\mathbf{u})) = \nabla_{\mathbf{z}} \mathcal{E}^*(\mathbf{z}(\mathbf{u})) = \mathbf{u}$

$$\mathcal{E}^{**}(\mathbf{u}) = \langle \mathbf{u}, \mathbf{z}(\mathbf{u}) \rangle - (\langle \mathbf{z}(\mathbf{u}), \mathbf{y}(\mathbf{z}(\mathbf{u})) \rangle - \mathcal{E}(\mathbf{y}(\mathbf{z}(\mathbf{u})))) = \mathcal{E}(\mathbf{u}).$$

Part 3. Now we add two infinitesimals $\boldsymbol{\tau}$ and $\boldsymbol{\delta}$ to \mathbf{z} and \mathbf{y} respectively s.t.

$$\nabla_{\mathbf{z}} \mathcal{E}^*(\mathbf{z} + \boldsymbol{\tau}) = \mathbf{y} + \boldsymbol{\delta}, \quad \nabla_{\mathbf{y}} \mathcal{E}(\mathbf{y} + \boldsymbol{\delta}) = \mathbf{z} + \boldsymbol{\tau}.$$

Then,

$$\nabla_{\mathbf{y}} \mathcal{E}(\mathbf{y} + \boldsymbol{\delta}) - \nabla_{\mathbf{y}} \mathcal{E}(\mathbf{y}) = \boldsymbol{\tau}, \quad \nabla_{\mathbf{z}} \mathcal{E}^*(\mathbf{z} + \boldsymbol{\tau}) - \nabla_{\mathbf{z}} \mathcal{E}^*(\mathbf{z}) = \boldsymbol{\delta}.$$

Since $\mathbf{H}_{\mathcal{E}}(\mathbf{y})$ measures the change of $\nabla_{\mathbf{y}} \mathcal{E}(\mathbf{y})$ when \mathbf{y} changes by an infinitesimal $\boldsymbol{\delta}$, then

$$\begin{aligned}
&\nabla_{\mathbf{y}} \mathcal{E}(\mathbf{y} + \boldsymbol{\delta}) - \nabla_{\mathbf{y}} \mathcal{E}(\mathbf{y}) \approx \mathbf{H}_{\mathcal{E}}(\mathbf{y}) \boldsymbol{\delta} \\
&\iff \mathbf{H}_{\mathcal{E}}^{-1}(\mathbf{y}) (\nabla_{\mathbf{y}} \mathcal{E}(\mathbf{y} + \boldsymbol{\delta}) - \nabla_{\mathbf{y}} \mathcal{E}(\mathbf{y})) \approx \boldsymbol{\delta} \\
&\iff \mathbf{H}_{\mathcal{E}}^{-1}(\mathbf{y}) \boldsymbol{\tau} \approx \boldsymbol{\delta} = \nabla \mathcal{E}^*(\mathbf{z} + \boldsymbol{\tau}) - \nabla \mathcal{E}^*(\mathbf{z}) \\
&\iff \mathbf{H}_{\mathcal{E}}^{-1}(\mathbf{y}) \boldsymbol{\tau} \approx \nabla \mathcal{E}^*(\mathbf{z} + \boldsymbol{\tau}) - \nabla \mathcal{E}^*(\mathbf{z})
\end{aligned} \tag{6}$$

Similarly,

$$\mathbf{H}_{\mathcal{E}^*}(\mathbf{z})\boldsymbol{\tau} \approx \nabla_{\mathbf{z}}\mathcal{E}^*(\mathbf{z} + \boldsymbol{\tau}) - \nabla_{\mathbf{z}}\mathcal{E}^*(\mathbf{z}) \quad (7)$$

Comparing (6) and (7), it is easy to see

$$\mathbf{H}_{\mathcal{E}^*}(\mathbf{z}) = \mathbf{H}_{\mathcal{E}}^{-1}(\mathbf{y}) \iff \mathbf{H}_{\mathcal{E}^*}(\nabla\mathcal{E}(\mathbf{y})) = \mathbf{H}_{\mathcal{E}}^{-1}(\mathbf{y}).$$

□

4 Newton's Method

4.1 Warm-up: Quadratic Optimization

First, let us play with a toy example, minimizing a quadratic function

$$\mathcal{E}(\mathbf{y}) = \frac{1}{2}\mathbf{y}^\top \mathbf{A}\mathbf{y} + \mathbf{b}^\top \mathbf{y} + c$$

where $\mathbf{A} \in \mathbb{R}^{n \times n}$ is positive definite. By setting the gradient w.r.t. \mathbf{y} to zero,

$$\nabla\mathcal{E}(\mathbf{y}) = \mathbf{A}\mathbf{y} + \mathbf{b} = \mathbf{0},$$

we obtain the global minimizer

$$\mathbf{y}^* = -\mathbf{A}^{-1}\mathbf{b}.$$

To make it more like gradient descent, let us start at some “guess” point \mathbf{y} and take a step $\boldsymbol{\delta}$ to move to the new point $\mathbf{y} + \boldsymbol{\delta}$. Then we try to minimize $\mathcal{E}(\mathbf{y} + \boldsymbol{\delta})$ by setting the gradient w.r.t. $\boldsymbol{\delta}$ to zero,

$$\begin{aligned} \nabla_{\boldsymbol{\delta}}\mathcal{E}(\mathbf{y} + \boldsymbol{\delta}) &= \mathbf{A}(\mathbf{y} + \boldsymbol{\delta}) = \mathbf{0} \\ \boldsymbol{\delta} &= -\mathbf{y} - \mathbf{A}^{-1}\mathbf{b} \\ \mathbf{y} + \boldsymbol{\delta} &= -\mathbf{A}^{-1}\mathbf{b} \end{aligned}$$

This gives us exactly global minimizer in just one step. However, the situation changes when the function is not quadratic anymore and thus we do not have constant Hessian. But taking a step which try to set the gradient to zero might still be a good idea.

4.2 K -stable Hessian

Next, consider a convex function $\mathcal{E} : \mathbb{R}^n \rightarrow \mathbb{R}$ whose Hessian is “nearly constant”. Recall the Hessian $\mathbf{H}_{\mathcal{E}}(\mathbf{y})$ aka $\nabla^2\mathcal{E}(\mathbf{y})$ at a point \mathbf{y} is just a matrix of pairwise 2nd order partial derivatives $\frac{\partial^2\mathcal{E}(\mathbf{y})}{\partial y_i\partial y_j}$. We say \mathcal{E} has k -stable Hessian if there exists a constant matrix \mathbf{A} s.t. for all \mathbf{y}

$$\mathbf{H}_{\mathcal{E}}(\mathbf{y}) \approx_K \mathbf{A} \iff \frac{1}{1+K}\mathbf{A} \preceq \mathbf{H}_{\mathcal{E}}(\mathbf{y}) \preceq (1+K)\mathbf{A}.$$

Note that we just require the existence of \mathbf{A} and do not assume we know \mathbf{A} . Then a natural question is to ask what convergence rate can be achieved if we take a gradient step “guided” by the

Hessian, which is called a “Newton step”. Such method is also known as the 2nd order method. Note that this is very similar to preconditioning.

Now, let us make our setting precise. We want minimize a convex function \mathcal{E} with k -stable Hessian $\mathbf{A} \succ \mathbf{0}$. And \mathbf{y}^* is a global minimizer of \mathcal{E} . Start from some initial point \mathbf{y}_0 . The update rule is

$$\mathbf{y}_{i+1} = \mathbf{y}_i - \alpha \cdot \mathbf{H}_{\mathcal{E}}^{-1}(\mathbf{y}_i) \nabla \mathcal{E}(\mathbf{y}_i),$$

where α is the step size and it will be decided later.

Theorem 4.1. $\mathcal{E}(\mathbf{y}_k) - \mathcal{E}(\mathbf{y}^*) \leq \epsilon (\mathcal{E}(\mathbf{y}_0) - \mathcal{E}(\mathbf{y}^*))$ when $k > (K + 1)^6 \log(1/\epsilon)$.

Proof. By Taylor’s theorem, there exists $\tilde{\mathbf{y}} \in [\mathbf{y}, \mathbf{y} + \delta]$ s.t.

$$\begin{aligned} \mathcal{E}(\mathbf{y} + \delta) &= \mathcal{E}(\mathbf{y}) + \nabla \mathcal{E}(\mathbf{y})^\top \delta + \frac{1}{2} \delta^\top \mathbf{H}_{\mathcal{E}}(\tilde{\mathbf{y}}) \delta \\ &\leq \underbrace{\mathcal{E}(\mathbf{y}) + \nabla \mathcal{E}(\mathbf{y})^\top \delta + \frac{(K+1)^2}{2} \delta^\top \mathbf{H}_{\mathcal{E}}(\mathbf{y}) \delta}_{=: f(\delta)} \end{aligned} \quad (8)$$

where the inequality comes from K -stability of Hessian,

$$\mathbf{H}_{\mathcal{E}}(\tilde{\mathbf{y}}) \preceq (1 + K) \mathbf{A} \preceq (1 + K)^2 \mathbf{H}_{\mathcal{E}}(\mathbf{y})$$

Observe that $f(\delta)$ is a convex quadratic function in δ . By setting minimizing it, or equivalently setting $\nabla_{\delta} f(\delta^*) = \mathbf{0}$, we get

$$\delta^* = -\frac{1}{(K+1)^2} \mathbf{H}_{\mathcal{E}}^{-1}(\mathbf{y}) \nabla_{\mathbf{y}} \mathcal{E}(\mathbf{y}) \quad (9)$$

Here, the step size α is equal to $(K+1)^{-2}$. Then, plugging (9) into (8),

$$\begin{aligned} \mathcal{E}(\mathbf{y} + \delta^*) &\leq \mathcal{E}(\mathbf{y}) - \frac{1}{2(K+1)^2} \nabla_{\mathbf{y}} \mathcal{E}(\mathbf{y})^\top \mathbf{H}_{\mathcal{E}}^{-1}(\mathbf{y}) \nabla_{\mathbf{y}} \mathcal{E}(\mathbf{y}) \\ &\leq \mathcal{E}(\mathbf{y}) - \frac{1}{2(K+1)^3} \nabla_{\mathbf{y}} \mathcal{E}(\mathbf{y})^\top \mathbf{A}^{-1} \nabla_{\mathbf{y}} \mathcal{E}(\mathbf{y}) \\ &\quad (\text{subtract } \mathcal{E}(\mathbf{y}^*) \text{ on both sides}) \\ \mathcal{E}(\mathbf{y} + \delta^*) - \mathcal{E}(\mathbf{y}^*) &\leq \mathcal{E}(\mathbf{y}) - \mathcal{E}(\mathbf{y}^*) - \frac{1}{2(K+1)^3} \underbrace{\nabla_{\mathbf{y}} \mathcal{E}(\mathbf{y})^\top \mathbf{A}^{-1} \nabla_{\mathbf{y}} \mathcal{E}(\mathbf{y})}_{=: \sigma} \end{aligned}$$

where the second inequality is due to K -stability of the inverse Hessian,

$$\frac{1}{1+K} \mathbf{A}^{-1} \preceq \mathbf{H}_{\mathcal{E}}(\mathbf{y})^{-1} \preceq (1+K) \mathbf{A}^{-1}.$$

Meanwhile, using Taylor’s theorem and K -stability, for some $\hat{\mathbf{y}}$ between \mathbf{y} and \mathbf{y}^* , and noting $\nabla \mathcal{E}(\mathbf{y}^*) = \mathbf{0}$, we have

$$\begin{aligned} \mathcal{E}(\mathbf{y}) &= \mathcal{E}(\mathbf{y}^*) + \nabla \mathcal{E}(\mathbf{y}^*)^\top (\mathbf{y} - \mathbf{y}^*) + \frac{1}{2} (\mathbf{y} - \mathbf{y}^*)^\top \mathbf{H}_{\mathcal{E}}(\hat{\mathbf{y}}) (\mathbf{y} - \mathbf{y}^*) \\ \mathcal{E}(\mathbf{y}) - \mathcal{E}(\mathbf{y}^*) &\leq \frac{(K+1)}{2} \underbrace{(\mathbf{y} - \mathbf{y}^*)^\top \mathbf{A} (\mathbf{y} - \mathbf{y}^*)}_{=: \gamma} \end{aligned}$$

Next, our task is reduced to comparing σ and γ . $\mathbf{y}_t := \mathbf{y}^* + t(\mathbf{y} - \mathbf{y}^*)$ ($t \in [0, 1]$) is a point on the segment connecting \mathbf{y}^* and \mathbf{y} . Since

$$\nabla \mathcal{E}(\mathbf{y}) = \nabla \mathcal{E}(\mathbf{y}) - \nabla \mathcal{E}(\mathbf{y}^*) = \int_0^1 H(\mathbf{y}_t)(\mathbf{y} - \mathbf{y}^*) dt,$$

then

$$\begin{aligned} (\mathbf{y} - \mathbf{y}^*)^\top \nabla \mathcal{E}(\mathbf{y}) &= \int_0^1 (\mathbf{y} - \mathbf{y}^*)^\top H(\mathbf{y}_t)(\mathbf{y} - \mathbf{y}^*) dt \\ &\geq \frac{1}{K+1} \int_0^1 (\mathbf{y} - \mathbf{y}^*)^\top \mathbf{A}(\mathbf{y} - \mathbf{y}^*) dt \\ &= \frac{\gamma}{K+1} \end{aligned} \tag{10}$$

On the other hand, define $\mathbf{z}_s = \nabla \mathcal{E}(\mathbf{y}^*) + s(\nabla \mathcal{E}(\mathbf{y}) - \nabla \mathcal{E}(\mathbf{y}^*))$ and then $d\mathbf{z}_s = \nabla \mathcal{E}(\mathbf{y}) ds$. Using *Theorem 3.3*, we have

$$\mathbf{y} - \mathbf{y}^* = \int_0^1 \mathbf{H}_{\mathcal{E}^*}(\mathbf{z}_s) \nabla \mathcal{E}(\mathbf{y}) ds.$$

Then,

$$\begin{aligned} \nabla \mathcal{E}(\mathbf{y})^\top (\mathbf{y} - \mathbf{y}^*) &= \int_0^1 \nabla \mathcal{E}(\mathbf{y})^\top \mathbf{H}_{\mathcal{E}^*}(\mathbf{z}_s) \nabla \mathcal{E}(\mathbf{y}) ds \\ &\leq (K+1) \int_0^1 \nabla \mathcal{E}(\mathbf{y})^\top \mathbf{A}^{-1} \nabla \mathcal{E}(\mathbf{y}) ds \\ &\leq (K+1)\sigma \end{aligned} \tag{11}$$

Combining (10) and (11) yields

$$\gamma \leq (K+1)^2 \sigma.$$

Therefore,

$$\mathcal{E}(\mathbf{y} + \boldsymbol{\delta}^*) - \mathcal{E}(\mathbf{y}^*) \leq (\mathcal{E}(\mathbf{y}) - \mathcal{E}(\mathbf{y}^*)) \left(1 - \frac{1}{(K+1)^6}\right).$$

□

Remark 4.2. The basic idea of relating σ and γ in the above proof is writing the same quantity, $\nabla \mathcal{E}(\mathbf{y})^\top (\mathbf{y} - \mathbf{y}^*)$, as two integrations along different lines. $(K+1)^6$ can be reduced to $(K+1)^2$ and even to $(K+1)$ with more care. In some settings, Newton's method converges in $\log \log(1/\epsilon)$ steps.

4.3 Linearly Constrained Newton's Method

Let us apply Newton's method to convex optimization problems with only linear constraints,

$$\begin{aligned} \min_{\mathbf{f} \in \mathbb{R}^m} \quad & \mathcal{E}(\mathbf{f}) \\ \text{s.t.} \quad & \mathbf{B}\mathbf{f} = \mathbf{d} \end{aligned}$$

where $\mathcal{E} : S \subseteq \mathbb{R}^m \rightarrow \mathbb{R}$ is a convex function and $\mathbf{B} \in \mathbb{R}^{n \times m}$. Wlog, let $\mathbf{d} = \mathbf{0}$, since otherwise we can equivalently deal the following problem with $\mathbf{B}\mathbf{f}_0 = \mathbf{d}$,

$$\begin{aligned} \min_{\boldsymbol{\rho} \in \mathbb{R}^m} \quad & \mathcal{E}(\mathbf{f}_0 + \boldsymbol{\rho}) \\ \text{s.t.} \quad & \mathbf{B}\boldsymbol{\rho} = \mathbf{0} \end{aligned}$$

It is useful to think the variable $\mathbf{f} \in \mathbb{R}^m$ as a flow in a graph. Define $C := \{\mathbf{f} : \mathbf{B}\mathbf{f} = \mathbf{0}\}$ which is essentially the kernel space of \mathbf{B} . C is also called the ‘‘cycle space’’ as it is the set of cycle flows when treating \mathbf{f} as flows. Restricting the domain of \mathcal{E} to the cycle space yields a new function $\hat{\mathcal{E}} : S \cap C \rightarrow \mathbb{R}$ s.t. $\hat{\mathcal{E}}(\mathbf{f}) = \mathcal{E}(\mathbf{f})$ for any $\mathbf{f} \in C$. How does $\nabla \hat{\mathcal{E}}$ look like compared to $\nabla \mathcal{E}$? Let Π_C be the orthogonal projection matrix onto C , meaning Π_C is symmetric and $\Pi_C \boldsymbol{\delta} = \boldsymbol{\delta}$ for any $\boldsymbol{\delta} \in C$. Given any $\mathbf{f} \in \mathbb{R}^m$, add to it an infinitesimal $\boldsymbol{\delta} \in C$, then

$$\begin{aligned} \mathcal{E}(\mathbf{f} + \boldsymbol{\delta}) &\approx \mathcal{E}(\mathbf{f}) + \langle \nabla \mathcal{E}(\mathbf{f}), \boldsymbol{\delta} \rangle \\ &= \mathcal{E}(\mathbf{f}) + \langle \nabla \mathcal{E}(\mathbf{f}), \Pi_C \boldsymbol{\delta} \rangle \\ &= \mathcal{E}(\mathbf{f}) + \langle \Pi_C \nabla \mathcal{E}(\mathbf{f}), \boldsymbol{\delta} \rangle \end{aligned}$$

This tells us the gradient of $\hat{\mathcal{E}}$ at a point $\mathbf{f} \in C$ is equal to the projection of gradient of $\nabla \mathcal{E}$ at \mathbf{f} onto the subspace C . Similarly,

$$\mathcal{E}(\mathbf{f} + \boldsymbol{\delta}) = \mathcal{E}(\mathbf{f}) + \langle \Pi_C \nabla \mathcal{E}(\mathbf{f}), \boldsymbol{\delta} \rangle + \frac{1}{2} \langle \boldsymbol{\delta}, \Pi_C \mathbf{H}_{\mathcal{E}}(\mathbf{f}) \Pi_C \boldsymbol{\delta} \rangle$$

Note that if \mathcal{E} has K -stable Hessian, then $\hat{\mathcal{E}}$ also has K -stable Hessian.

What is a Newton step $\boldsymbol{\delta}^*$ in a linearly constrained optimization problem? $\boldsymbol{\delta}^*$ should be a minimizer of

$$\begin{aligned} \min_{\substack{\boldsymbol{\delta} \in \mathbb{R}^m \\ \mathbf{B}\boldsymbol{\delta} = \mathbf{0}}} \quad & \underbrace{\langle \nabla \mathcal{E}(\mathbf{f}), \boldsymbol{\delta} \rangle}_{=: \mathbf{g}} + \frac{1}{2} \underbrace{\langle \boldsymbol{\delta}, \mathbf{H}_{\mathcal{E}}(\mathbf{f}) \boldsymbol{\delta} \rangle}_{=: \mathbf{H}} \\ & \text{(Lagrange duality)} \\ \iff \max_{\mathbf{x} \in \mathbb{R}^n} \min_{\boldsymbol{\delta} \in \mathbb{R}^m} \quad & \underbrace{\langle \mathbf{g}, \boldsymbol{\delta} \rangle + \frac{1}{2} \langle \boldsymbol{\delta}, \mathbf{H}\boldsymbol{\delta} \rangle - \mathbf{x}^\top \mathbf{B}\boldsymbol{\delta}}_{\text{Lagrangian } L(\boldsymbol{\delta}, \mathbf{x})} \end{aligned} \tag{12}$$

Applying the KKT optimality conditions, one has

$$\begin{aligned} \mathbf{B}\boldsymbol{\delta} &= \mathbf{0}, \\ \nabla_{\boldsymbol{\delta}} L(\boldsymbol{\delta}, \mathbf{x}) &= \mathbf{g} + \mathbf{H}\boldsymbol{\delta} - \mathbf{B}^\top \mathbf{x} = \mathbf{0}, \end{aligned}$$

from which we get

$$\begin{aligned} \boldsymbol{\delta} + \mathbf{H}^{-1} \mathbf{g} &= \mathbf{H}^{-1} \mathbf{B}^\top \mathbf{x} \\ \underbrace{\mathbf{B}\boldsymbol{\delta}}_{=0} + \mathbf{B}\mathbf{H}^{-1} \mathbf{g} &= \mathbf{B}\mathbf{H}^{-1} \mathbf{B}^\top \mathbf{x} \\ \mathbf{B}\mathbf{H}^{-1} \mathbf{g} &= \underbrace{\mathbf{B}\mathbf{H}^{-1} \mathbf{B}^\top}_{=: \mathbf{L}} \mathbf{x} \end{aligned}$$

Finally, the solutions to (12) are

$$\begin{cases} \mathbf{x}^* &= \mathbf{L}^{-1} \mathbf{B} \mathbf{H}^{-1} \mathbf{g} \\ \boldsymbol{\delta}^* &= -\mathbf{H}^{-1} \mathbf{g} + \mathbf{H}^{-1} \mathbf{B}^\top \mathbf{x}^* \end{cases}$$

It is easy to verify that $\mathbf{B} \boldsymbol{\delta}^* = \mathbf{0}$. Thus, our update rule is $\mathbf{f}_{i+1} = \mathbf{f}_i + \boldsymbol{\delta}^*$. And we have the following convergence result.

Theorem 4.3. $\hat{\mathcal{E}}(\mathbf{f}_k) - \hat{\mathcal{E}}(\mathbf{f}^*) \leq \epsilon \cdot \left(\hat{\mathcal{E}}(\mathbf{f}_0) - \hat{\mathcal{E}}(\mathbf{f}^*) \right)$ when $k > (K + 1)^6 \log(1/\epsilon)$.

Remark 4.4. Note if $\mathcal{E}(\mathbf{f}) = \sum_{i=1}^m \mathcal{E}_i(\mathbf{f}(i))$, then $\mathbf{H}_{\mathcal{E}}(\mathbf{f})$ is diagonal. Thus, $\mathbf{L} = \mathbf{B} \mathbf{H}^{-1} \mathbf{B}^\top$ is indeed a Laplacian provided that \mathbf{B} is an incidence matrix. Therefore, the linear equations we need to solve to apply Newton's method in a network flow setting are Laplacians, which means we can solve them very quickly.