

A Tutorial on Gradient Descend

Richard Xu

August 29, 2022

1 Implicit bias of gradient descend

This section explains [1]. The big picture here is to show the gradient $\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) \neq \mathbf{0}$ (section of 1.3.2), the loss function $\mathcal{L}(\mathbf{w})$ will continue to decrease using gradient descent. This makes $\|\mathbf{w}(t)\| \rightarrow \infty$ as $t \rightarrow \infty$. As a result, the weights of the few dominant linear combination terms correspond to the weights associated with the support vectors.

1.1 classifier without max-margin

looking at support vector machine term below:

$$\begin{aligned} \min \left(\frac{1}{2} \|\mathbf{w}\|^2 \right) \\ \text{subject to: } 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + w_0) \leq 0 \quad \forall i \end{aligned} \quad (1)$$

If we were not trying to solve a max-margin problem: if we were just trying to express the problem as a linear classifier. Then, the objective (for a single \mathbf{x}_i, y_i pair can be written as):

$$y_i(\mathbf{w}^\top \mathbf{x}_i + w_0) > 0 \quad (2)$$

to make things even simpler, drop the w_0 :

$$y_i(\mathbf{w}^\top \mathbf{x}_i) > 0 \quad (3)$$

1.1.1 smooth loss

smooth loss function used to penalize incorrect classification, for example:

$$\begin{aligned} \ell(u) &= \exp^{-u} \\ \implies \ell(\mathbf{w}^\top \mathbf{x}_i y_i) &= \exp(-\mathbf{w}^\top \mathbf{x}_i y_i) \end{aligned} \quad (4)$$

in words, we must “push” value of $\mathbf{w}^\top \mathbf{x}_i y_i$ to be large +ve value (for correctly classified data/label pairs) when smooth loss function is assigned to

1.2 use gradient descend

when gradient descend is used to minimize the objective below (note analytical solution available for svm):

$$\begin{aligned}
& \min \mathcal{L}(\mathbf{w}) \\
& = \min \sum_{i=1}^n \ell(\mathbf{w}^\top \mathbf{x}_i y_i) \\
& = \min \sum_{i=1}^n \ell(\mathbf{w}^\top \tilde{\mathbf{x}}_i) \quad \text{let } \tilde{\mathbf{x}}_i = \mathbf{x}_i y_i
\end{aligned} \tag{5}$$

1.2.1 gradient for generic loss \mathcal{L}

$$\begin{aligned}
\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) &= \nabla_{\mathbf{w}} \sum_{i=1}^n \ell(\mathbf{w}^\top \tilde{\mathbf{x}}_i) \\
&= \sum_{i=1}^n \ell'(\mathbf{w}^\top \tilde{\mathbf{x}}_i) \tilde{\mathbf{x}}_i
\end{aligned} \tag{6}$$

substitute into gradient descend:

$$\begin{aligned}
\mathbf{w}(t+1) &= \mathbf{w}(t) - \eta \nabla_{\mathbf{w}(t)} \mathcal{L}(\mathbf{w}(t)) \\
&= \mathbf{w}(t) - \eta \sum_{i=1}^n \ell'(\mathbf{w}(t)^\top \tilde{\mathbf{x}}_i) \tilde{\mathbf{x}}_i
\end{aligned} \tag{7}$$

we are interested in the behavior of $\mathbf{w}(t) \rightarrow \infty$

1.3 magnitude: $\|\mathbf{w}(t)\| \rightarrow \infty$

1.3.1 no finite critical points $\nabla_{\mathbf{w}(t)} \mathcal{L}(\mathbf{w}(t)) = 0$

It's difficult to show from the gradient directly why the expression $\sum_{i=1}^n \ell'(\mathbf{w}(t)^\top \tilde{\mathbf{x}}_i) \tilde{\mathbf{x}}_i$ never reach 0, i.e.,

$$\text{to show why } \lim_{t \rightarrow \infty} \sum_{i=1}^n \ell'(\mathbf{w}(t)^\top \tilde{\mathbf{x}}_i) \tilde{\mathbf{x}}_i \neq 0 \tag{8}$$

Note that people may be confused to think if we let $\ell(u) = \exp^{-u}$, then $\ell'(u) \neq 0$ anyway. right? However, since we have a sum and not just a term. Making the gradient zero may still seems "possible". To illustrate, when we let $n = 2$, we may obtain a situation where:

$$\ell'(\mathbf{w}^\top \tilde{\mathbf{x}}_1) \tilde{\mathbf{x}}_1 = -\ell'(\mathbf{w}^\top \tilde{\mathbf{x}}_2) \tilde{\mathbf{x}}_2 \quad \text{for some } \mathbf{w} \tag{9}$$

1.3.2 show $\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}(t))$ won't be a zero vector

Let's assume $\exists \mathbf{w}^* \neq \mathbf{0}$ making data separable (if data is separable). looking at the following expression:

$$\begin{aligned} \mathbf{w}^{*\top} \eta \nabla_{\mathbf{w}(t)} \mathcal{L}(\mathbf{w}(t)) &= \mathbf{w}^{*\top} \sum_{i=1}^n \ell'(\mathbf{w}(t)^\top \tilde{\mathbf{x}}_i) \tilde{\mathbf{x}}_i \\ &= \sum_{i=1}^n \underbrace{\ell'(\mathbf{w}(t)^\top \tilde{\mathbf{x}}_i)}_{<0} \underbrace{\tilde{\mathbf{x}}_i^\top \mathbf{w}^*}_{>0} \end{aligned} \quad (10)$$

Obviously, since:

$$\begin{aligned} \ell'(\mathbf{w}(t)^\top \tilde{\mathbf{x}}_i) \tilde{\mathbf{x}}_i^\top \mathbf{w}^* &< 0 \text{ and } \mathbf{w}^* \neq \mathbf{0} \\ \implies \ell'(\mathbf{w}(t)^\top \tilde{\mathbf{x}}_i) \tilde{\mathbf{x}}_i &\neq \mathbf{0} \quad \forall i \\ \implies \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}(t)) &\neq \mathbf{0} \end{aligned} \quad (11)$$

Explain each two terms:

1. $\mathbf{w}^{*\top} \tilde{\mathbf{x}}_i > 0 \quad \forall i$ if all data are all correctly classified/linearly separable:

$$y_i(\mathbf{w}^{*\top} \mathbf{x}_i) > 0 \quad (12)$$

note that up to here, we made **no** reference with max-margin

2. $\ell'(\cdot) < 0$ as long as we choose a monotonically decreasing ℓ which means its gradient < 0
3. also note that in here, we merely assumed $\exists \mathbf{w}^*$. Don't get confused, it is not where $\mathbf{w}(t)$ converges to!
4. also note if it's possible for $\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}(t)) = \mathbf{0}$, it means the gradient descend will not run indefinitely.

1.3.3 why $\|\mathbf{w}(t)\| \rightarrow \infty$?

We know gradient descend on a smooth loss will converge to a minimum. This will be illustrated in the β -smooth section. Since ℓ is a smooth function, so is $\mathcal{L}(\mathbf{w}(t)) = \sum_{i=1}^n \ell(\mathbf{w}(t)^\top \tilde{\mathbf{x}}_i)$:

$$\begin{aligned}
\|\nabla \mathcal{L}(x) - \nabla \mathcal{L}(y)\| &= \left\| \frac{1}{n} \sum_i \nabla \ell_i(x) - \frac{1}{n} \sum_i \nabla \ell_i(y) \right\| \\
&= \frac{1}{n} \left\| \sum_i (\nabla \ell_i(x) - \nabla \ell_i(y)) \right\| \\
&\leq \frac{1}{n} \sum_i \left\| \nabla \ell_i(x) - \nabla \ell_i(y) \right\| \quad \text{triangle inequality} \quad (13) \\
&\leq \frac{1}{n} \sum_i (\beta_i \|x - y\|) \\
&= \left(\frac{1}{n} \sum_i \beta_i \right) \|x - y\|
\end{aligned}$$

However, the above says there is no critical points. Putting above two arguments together, and look at the objective $\sum_{i=1}^n \ell(\mathbf{w}^\top \tilde{\mathbf{x}}_i)$, we can see that, since the gradient descend algorithm continues to run (and the loss will continuously becoming smaller):

$$\left(\mathcal{L}(\mathbf{w}(t)) = \sum_{i=1}^n \ell(\mathbf{w}(t)^\top \tilde{\mathbf{x}}_i) \right) \rightarrow 0 \implies \mathbf{w}(t)^\top \tilde{\mathbf{x}}_i \rightarrow \infty \quad \text{think } \exp(-u) \quad (14)$$

Since $\tilde{\mathbf{x}}_i$ is fixed, then $\|\mathbf{w}(t)\| \rightarrow \infty$. Note that this is why we need to show there is **no** critical points first.

The norm is needed as $y_i \in \{1, -1\}$, it means:

$$\begin{aligned}
&\lim_{t \rightarrow \infty} \|\mathbf{w}(t)\| = \infty \\
&\text{or equivalently } \|\mathbf{w}(t)\| \rightarrow \infty \quad (15)
\end{aligned}$$

1.4 what about direction of $\mathbf{w}(t)$?

To characterize direction, we look at normalized $\lim_{t \rightarrow \infty} \frac{\mathbf{w}(t)}{\|\mathbf{w}(t)\|}$

Theorem 1 *under assumption as $t \rightarrow \infty$, Gradient descend behaves as:*

$$\mathbf{w}(t) \approx \frac{\mathbf{w}_{\text{svm}}}{\|\mathbf{w}_{\text{svm}}\|} \quad (16)$$

1.4.1 explanation

when $\mathbf{w}(t) \rightarrow \infty$, it has the same direction of the SVM solution, i.e., its normalized version $\frac{\mathbf{w}(t)}{\|\mathbf{w}(t)\|}$ becomes that of the \mathbf{w}_{svm}

\mathbf{w}_{svm} gives max-margin classifier which has better generalization!

1.5 proof of theorem

consider exponential loss $\mathcal{L}(u) = \exp(-u)$, gradient descend in asymptotic regime in shown in Eq.(14):

$$\mathbf{w}(t)^\top \tilde{\mathbf{x}}_i \rightarrow \infty \quad \forall i \quad (17)$$

1.5.1 what is asymptotic “simplification” convergence?

The definition of the notation $a_n \rightarrow b_n$ is designed to mean that $a_n \approx b_n$ for large n , where the fit gets better and better as n gets larger, for example:

$$\lim_{x \rightarrow \infty} x^2 + x + 1 = x^2 \quad (18)$$

and,

$$\begin{aligned} u(x+h) &= u(x) + u'(x)h + \frac{u''(x)}{2}h^2 + \dots \\ u(x+h) - u(x) + u'(x)h &= \frac{u''(x)}{2}h^2 + \dots \\ \left| \frac{u(x+h) - u(x)}{h} - u'(x) \right| &= \left| \frac{u''(x)}{2}h + \frac{u'''(x)}{3!}h^2 \dots \right| \quad \text{divided by } h \\ \Rightarrow \lim_{h \rightarrow 0} \left| \frac{u(x+h) - u(x)}{h} - u'(x) \right| &= \left| \frac{u''(x)}{2}h \right| \end{aligned} \quad (19)$$

1.5.2 asymptotic convergence of $\frac{\mathbf{w}(t)}{\|\mathbf{w}(t)\|}$

we express $\mathbf{w}(t)$ as a linear function in terms of \mathbf{w}_∞ (the remaining work is to find what \mathbf{w}_∞ is):

$$\mathbf{w}(t) = \underbrace{m(t)}_{\text{magnitude}} \mathbf{w}_\infty + \underbrace{\mathbf{b}(t)}_{\text{residual}} \quad (20)$$

assume $\exists \mathbf{w}_\infty$ (which is a unit vector), the limit of the normalization $\frac{\mathbf{w}(t)}{\|\mathbf{w}(t)\|} \rightarrow \mathbf{w}_\infty$, under assumptions of both already stated, and new ones:

1. $\lim_{t \rightarrow \infty} \frac{\mathbf{b}(t)}{m(t)} = 0$ as $\|\mathbf{b}(t)\|$ is relatively smaller compare with $\|\mathbf{w}(t)\|$, as $t \rightarrow \infty$
2. $m(t) \rightarrow \infty$ makes sense as $\|\mathbf{w}(t)\| \rightarrow \infty$

since $m(t)$ is the magnitude, then $m(t) \geq 0$. Looking at the gradient again:

$$\begin{aligned}
\nabla_{\mathbf{w}(t)} \mathcal{L}(\mathbf{w}(t)) &= \sum_{i=1}^n \ell'(\mathbf{w}(t)^\top \tilde{\mathbf{x}}_i) \tilde{\mathbf{x}}_i \\
&= - \sum_{i=1}^n \exp(-\mathbf{w}(t)^\top \tilde{\mathbf{x}}_i) \tilde{\mathbf{x}}_i \quad \text{substitute } \ell'(u) = -\exp(-u) \\
\Rightarrow -\nabla_{\mathbf{w}(t)} \mathcal{L}(\mathbf{w}(t)) &= \sum_{i=1}^n \exp(-\mathbf{w}(t)^\top \tilde{\mathbf{x}}_i) \tilde{\mathbf{x}}_i \\
&= \sum_{i=1}^n \exp\left(-\left(m(t)\mathbf{w}_\infty + \mathbf{b}(t)\right)^\top \tilde{\mathbf{x}}_i\right) \tilde{\mathbf{x}}_i \quad \text{substitute } \mathbf{w}(t) = m(t)\mathbf{w}_\infty + \mathbf{b}(t) \\
&= \sum_{i=1}^n \exp^{-m(t)\mathbf{w}_\infty^\top \tilde{\mathbf{x}}_i} \tilde{\mathbf{x}}_i \times \exp^{-\mathbf{b}(t)^\top \tilde{\mathbf{x}}_i} \tilde{\mathbf{x}}_i \\
&\approx \sum_{i=1}^n \exp^{-m(t)\mathbf{w}_\infty^\top \tilde{\mathbf{x}}_i} \tilde{\mathbf{x}}_i \quad \because \lim_{t \rightarrow \infty} \frac{\mathbf{b}(t)}{m(t)} = 0 \\
&= \sum_{i=1}^n \underbrace{\exp\left(-m(t)\mathbf{w}_\infty^\top \tilde{\mathbf{x}}_i\right)}_{\alpha_i} \tilde{\mathbf{x}}_i
\end{aligned} \tag{21}$$

so gradient step would be some non-negative linear combination of $\tilde{\mathbf{x}}_i$, i.e.,:

$$-\nabla_{\mathbf{w}(t)} \mathcal{L}(\mathbf{w}(t)) = \sum_{i=1}^n \alpha_i \tilde{\mathbf{x}}_i \tag{22}$$

1.5.3 dominate terms

assumes \mathbf{w}_∞ classifies the linearly separable data correctly, then:

$$\mathbf{w}_\infty^\top \tilde{\mathbf{x}}_i > 0 \tag{23}$$

since $m(t) \rightarrow \infty$, we have only a few dominate terms in $\{\tilde{\mathbf{x}}_i\}$ (multiply by ∞ makes them dominate!). since these $\tilde{\mathbf{x}}_i$ are closest to (and on the) decision boundary, then they are precisely support vectors! So the set is support vector set s.v.!

Note that if multiple $\tilde{\mathbf{x}}_i$ and $\tilde{\mathbf{x}}_j$ are the closest, i.e., $\mathbf{w}_\infty^\top \tilde{\mathbf{x}}_i = \mathbf{w}_\infty^\top \tilde{\mathbf{x}}_j$, then, they both are part of the support vector set!

$$\begin{aligned}
-\nabla_{\mathbf{w}(t)} \mathcal{L}(\mathbf{w}(t)) &\approx \sum_{\tilde{\mathbf{x}}_i \in \text{s. v.}} \alpha_i \tilde{\mathbf{x}}_i \\
&= \sum_{\mathbf{x}_i \in \text{s. v.}} \alpha_i \mathbf{x}_i y_i
\end{aligned} \tag{24}$$

As each of the gradient step is a linear combination of $x_i \in \text{s.v.}$, then, so is \mathbf{w}_∞ , i.e.,

$$\mathbf{w}_\infty = \sum_{\mathbf{x}_i \in \text{s.v.}} \alpha'_i \tilde{\mathbf{x}}_i \quad \text{for some } \alpha'_i \neq \alpha_i \tag{25}$$

since $\|\mathbf{w}(t)\| \rightarrow \infty$, then the initial $\mathbf{w}(0)$ value won't matter any more. There is one remaining issue though: $\mathbf{w}_\infty^\top \tilde{\mathbf{x}}_i \neq 1 \quad \forall \mathbf{x}_i \in \text{s.v.}$ so look at the next section:

1.5.4 from \mathbf{w}_∞ to obtain \mathbf{w}_{svm}

lastly, we need to scale \mathbf{w}_∞ to become \mathbf{w}_{svm} . let's see what if we perform $\frac{\mathbf{w}_\infty}{\text{some constant}}$. Now let's have $\tilde{\mathbf{x}}_{\text{s.v.}}$ such that:

$$\mathbf{w}_\infty^\top \tilde{\mathbf{x}}_{\text{s.v.}} = \min_i \{\mathbf{w}_\infty^\top \tilde{\mathbf{x}}_i\} \quad (26)$$

although the picking of the "some constant" is arbitrary, but we pick $\min_i \{\mathbf{w}_\infty^\top \tilde{\mathbf{x}}_i\}$ to reflect the SVM solution:

$$\begin{aligned} \hat{\mathbf{w}} &= \frac{\mathbf{w}_\infty}{\text{some constant}} \\ &= \frac{\mathbf{w}_\infty}{\mathbf{w}_\infty^\top \tilde{\mathbf{x}}_{\text{s.v.}}} \end{aligned} \quad (27)$$

note that $\|\mathbf{w}_\infty\| = 1$, but $\|\hat{\mathbf{w}}\| \neq 1$! By this process, it scales $\hat{\mathbf{w}}$ such that when applying to $\tilde{\mathbf{x}}_{\text{s.v.}}$:

$$\begin{aligned} \hat{\mathbf{w}}^\top \tilde{\mathbf{x}}_{\text{s.v.}} &= \frac{\mathbf{w}_\infty^\top}{\mathbf{w}_\infty^\top \tilde{\mathbf{x}}_{\text{s.v.}}} \tilde{\mathbf{x}}_{\text{s.v.}} \\ &= 1 \end{aligned} \quad (28)$$

and when it applies to other $\tilde{\mathbf{x}} \notin \{\tilde{\mathbf{x}}_{\text{s.v.}}\}$:

$$\begin{aligned} \hat{\mathbf{w}}^\top \tilde{\mathbf{x}} &= \frac{\mathbf{w}_\infty^\top}{\mathbf{w}_\infty^\top \tilde{\mathbf{x}}_{\text{s.v.}}} \tilde{\mathbf{x}} \\ &> 1 \end{aligned} \quad (29)$$

Does $\hat{\mathbf{w}}$ look familiar? Remember KKT condition is:

$$\hat{\mathbf{w}} = \sum_{i=1}^N \lambda_i \tilde{\mathbf{x}}_i \quad (30)$$

with complementary duality:

$$\begin{cases} \lambda_i > 0 & \hat{\mathbf{w}}^\top \tilde{\mathbf{x}}_i = 1 & \text{support vectors} \\ \lambda_i = 0 & \hat{\mathbf{w}}^\top \tilde{\mathbf{x}}_i > 1 & \text{non support vector} \end{cases} \quad (31)$$

compare with equation in SVM section, $\hat{\mathbf{w}} = \mathbf{w}_{\text{svm}}$

Since we already prove $\hat{\mathbf{w}}$ is proportional to \mathbf{w}_∞ . Therefore, \mathbf{w}_∞ is the SVM solution up to some constant!

2 Fenchel dual function

$$f^*(\mathbf{y}) = \sup_{\mathbf{x}} [\mathbf{x}^\top \mathbf{y} - f(\mathbf{x})] \quad (32)$$

2.1 property of Fenchel dual

2.1.1 visualization

visualization is achieved similarly to what was done previously in the general dual function section: Consider $f^*(\mathbf{y})$ is a function on the “gradient space”, except we now have parameter:

$$\lambda \rightarrow \mathbf{y} \quad (33)$$

Just like the general dual function section, we can visualize by generating $f^*(\mathbf{y})$ from maximization of finite lines (in \mathbf{y}) defined by a finite set $\{\mathbf{x}\}$, and \mathbf{x} and $f(\mathbf{x})$ are treated like “constant line parameters”. the alternative way to consider this is to rewrite Eq.(32) as:

$$f^*(\mathbf{y}) = \sup_{\mathbf{x}} [-f(\mathbf{x}) + \mathbf{y}^\top \mathbf{x}] \quad (34)$$

where the dual “constraint” $g(\mathbf{x}) \equiv \mathbf{y}^\top \mathbf{x}$. Note $f^*(\mathbf{y})$ is always convex.

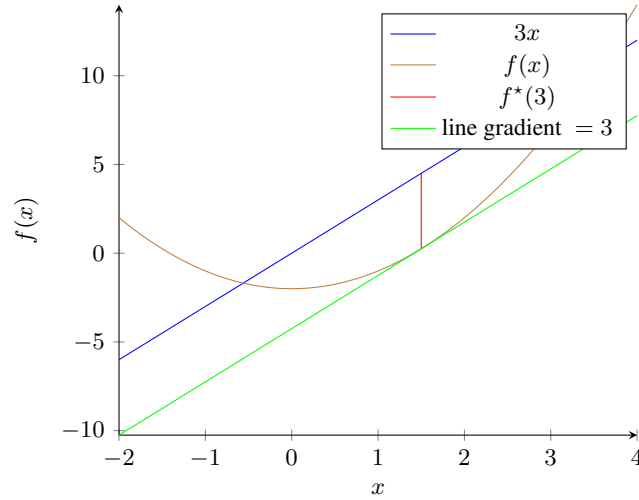


Figure 1: simple demonstration of Fenchel dual function

imagine we picked a particular value of $\mathbf{y} = 3$ in order to compute the value of $f^*(3)$, then we worked out:

$$\begin{aligned} f^*(3) &= \sup_x [3x - f(x)] \\ \hat{x}_3 &= \arg \max_x [3x - f(x)] \end{aligned} \quad (35)$$

which is the **red** line segment, i.e., the line segment/difference between the line $3x$ (gradient 3, passing through origin), i.e., the **blue** line and the function $f(x)$, i.e., the **brown** line.

2.1.2 important observation

introducing arg max variable \hat{x}_y , such that by given y :

$$\begin{aligned}\hat{x}_y &= \arg \max_x [\mathbf{y}^\top \mathbf{x} - f(\mathbf{x})] \\ \implies \nabla_{\mathbf{x}} (\hat{x}_y^\top \mathbf{y} - f(\hat{x}_y)) &= 0 \\ \implies \nabla_{\mathbf{x}} f(\hat{x}_y) &= \mathbf{y}\end{aligned}\tag{36}$$

It says that given y , maximum of the line segment length between $\mathbf{y}^\top \mathbf{x}$ and $f(\mathbf{x})$ occurs at particular \hat{x}_y , where its gradient $\nabla_{\mathbf{x}} f(\hat{x}_y) = \mathbf{y}$. In our case, we picked $y = 3$, therefore $\nabla_x f(\hat{x}_3) = 3$.

Visually, we see two parallel lines in **blue** and **green** as both have gradient $y = 3$. So far, we express everything in term of gradients. However, they will be replaced by sub-gradients.

2.1.3 Fenchel's inequality

for any \mathbf{x} and y :

$$f(\mathbf{x}) + f^*(\mathbf{y}) \geq \mathbf{x}^\top \mathbf{y}\tag{37}$$

the reason is because:

$$\begin{aligned}f^*(\mathbf{y}) &= \sup_{\mathbf{x}} [\mathbf{x}^\top \mathbf{y} - f(\mathbf{x})] \\ &\geq \mathbf{x}^\top \mathbf{y} - f(\mathbf{x}) \\ \implies f(\mathbf{x}) + f^*(\mathbf{y}) &\geq \mathbf{x}^\top \mathbf{y}\end{aligned}\tag{38}$$

2.1.4 example

$$\begin{aligned}f(\mathbf{x}) &= \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} + \mathbf{q}^\top \mathbf{x} + c \\ f^*(\mathbf{y}) &= \sup_{\mathbf{x}} \left(\mathbf{y}^\top \mathbf{x} - \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{q}^\top \mathbf{x} - c \right)\end{aligned}\tag{39}$$

to find optimal \mathbf{x} :

$$\begin{aligned}\hat{x}_y - \mathbf{Q} \mathbf{x} - \mathbf{q} &= \mathbf{0} \\ \hat{x}_y &= \mathbf{Q}^{-1}(\mathbf{y} - \mathbf{q})\end{aligned}\tag{40}$$

substitute it back to $f^*(\mathbf{y})$:

$$\begin{aligned}
f^*(\mathbf{y}) &= \mathbf{y}^\top (\mathbf{Q}^{-1}(\mathbf{y} - \mathbf{q})) - \frac{1}{2} (\mathbf{Q}^{-1}(\mathbf{y} - \mathbf{q}))^\top \mathbf{Q} (\mathbf{Q}^{-1}(\mathbf{y} - \mathbf{q})) - \mathbf{q}^\top (\mathbf{Q}^{-1}(\mathbf{y} - \mathbf{q})) - c \\
&= (\mathbf{y} - \mathbf{q})^\top (\mathbf{Q}^{-1}(\mathbf{y} - \mathbf{q})) - \frac{1}{2} (\mathbf{Q}^{-1}(\mathbf{y} - \mathbf{q}))^\top \mathbf{Q} (\mathbf{Q}^{-1}(\mathbf{y} - \mathbf{q})) - c \\
&= (\mathbf{y} - \mathbf{q})^\top \mathbf{Q}^{-1}(\mathbf{y} - \mathbf{q}) - \frac{1}{2} (\mathbf{y} - \mathbf{q})^\top \mathbf{Q}^{-\top} \mathbf{Q} \mathbf{Q}^{-1}(\mathbf{y} - \mathbf{q}) - c \\
&= (\mathbf{y} - \mathbf{q})^\top \mathbf{Q}^{-1}(\mathbf{y} - \mathbf{q}) - \frac{1}{2} (\mathbf{y} - \mathbf{q})^\top \mathbf{Q}^{-1}(\mathbf{y} - \mathbf{q}) - c \quad \because \mathbf{Q}^\top = \mathbf{Q} \\
&= \frac{1}{2} (\mathbf{y} - \mathbf{q})^\top \mathbf{Q}^{-1}(\mathbf{y} - \mathbf{q}) - c
\end{aligned} \tag{41}$$

2.1.5 conjugate of conjugate function in general

In general, for **any** arbitrary function $f(x)$, conjugate of conjugate function is no greater than the original function:

$$f^{**}(\mathbf{x}) \leq f(\mathbf{x}) \tag{42}$$

this can be easily proven:

$$\begin{aligned}
f^{**}(\mathbf{x}) &= \sup_{\mathbf{y}} [\mathbf{x}^\top \mathbf{y} - f^*(\mathbf{y})] \\
&= \sup_{\mathbf{y}} [\mathbf{x}^\top \mathbf{y} - \sup_{\mathbf{z}} [\mathbf{y}^\top \mathbf{z} - f(\mathbf{z})]] \quad \text{need to use variable } \mathbf{z} \\
&= \sup_{\mathbf{y}} [\mathbf{x}^\top \mathbf{y} + \inf_{\mathbf{z}} [-\mathbf{y}^\top \mathbf{z} + f(\mathbf{z})]] \quad -\sup\{\mathbf{x}\} \Leftrightarrow +\inf\{-\mathbf{x}\} \\
&= \sup_{\mathbf{y}} [\inf_{\mathbf{z}} [\mathbf{y}^\top (\mathbf{z} - \mathbf{x}) + f(\mathbf{z})]] \\
&\leq \inf_{\mathbf{z}} [\sup_{\mathbf{y}} [\mathbf{y}^\top (\mathbf{z} - \mathbf{x}) + f(\mathbf{z})]] \\
&= f(\mathbf{x})
\end{aligned} \tag{43}$$

the last line can be seen as \mathbf{y} is not bounded, so it can always be chosen that the inner term $\sup_{\mathbf{y}} [\mathbf{y}^\top (\mathbf{z} - \mathbf{x}) + f(\mathbf{z})] \rightarrow \infty$, therefore, we can prevent this by having $\mathbf{z} = \mathbf{x}$. Note that this is used often in the context of min max or inf sup.

2.1.6 conjugate of conjugate function when $f(\mathbf{x})$ is convex

however, when $f(\mathbf{x})$ is convex:

$$f^{**}(\mathbf{x}) = f(\mathbf{x}) \tag{44}$$

we can exploit the above property:

$$\begin{aligned}
f^*(\mathbf{y}) &= \max_{\mathbf{x}} [\mathbf{x}^\top \mathbf{y} - f(\mathbf{x})] \quad \forall \mathbf{y} \quad \text{by definition} \\
\Rightarrow f(\mathbf{x}) &= f^{**}(\mathbf{x}) = \max_{\mathbf{y}} [\mathbf{x}^\top \mathbf{y} - f^*(\mathbf{y})] \quad \forall \mathbf{x} \quad \text{by } f^{**}(\mathbf{x}) = f(\mathbf{x})
\end{aligned} \tag{45}$$

2.1.7 derivatives of conjugate function

when $f(\mathbf{x})$ is convex and differentiable, then Eq.(36) can be extended to both ways:

$$\begin{aligned} f(\mathbf{x}) + f^*(\mathbf{y}) &= \mathbf{x}^\top \mathbf{y} \quad \forall \mathbf{x}, \mathbf{y} \\ \implies \nabla_{\mathbf{x}} f(\mathbf{x}) &= \mathbf{y} \quad \text{and} \quad \nabla_{\mathbf{y}} f^*(\mathbf{y}) = \mathbf{x} \end{aligned} \quad (46)$$

note that we make the two equations of the last line having the same \mathbf{x} and \mathbf{y} , we call them $\hat{\mathbf{x}}_{\mathbf{y}}$ and $\hat{\mathbf{y}}_{\mathbf{x}}$:

$$\nabla_{\mathbf{x}} f(\hat{\mathbf{x}}_{\mathbf{y}}) = \hat{\mathbf{y}}_{\mathbf{x}} \quad \text{and} \quad \nabla_{\mathbf{y}} f^*(\hat{\mathbf{y}}_{\mathbf{x}}) = \hat{\mathbf{x}}_{\mathbf{y}} \quad (47)$$

$$\nabla_{\mathbf{x}} f(\underbrace{\nabla_{\mathbf{y}} f^*(\hat{\mathbf{y}}_{\mathbf{x}})}_{\hat{\mathbf{x}}_{\mathbf{y}}}) = \hat{\mathbf{y}}_{\mathbf{x}} \quad (48)$$

therefore, for a corresponding pair $(\mathbf{x}_{\mathbf{y}}, \mathbf{y}_{\mathbf{x}})$, obtained from computing the conjugate dual:

$$\nabla_{\mathbf{x}} f(\hat{\mathbf{x}}_{\mathbf{y}}) = \hat{\mathbf{y}}_{\mathbf{x}} \quad \iff \quad \nabla_{\mathbf{y}} f^*(\hat{\mathbf{y}}_{\mathbf{x}}) = \hat{\mathbf{x}}_{\mathbf{y}} \quad (49)$$

This property says that the function argument $\hat{\mathbf{x}}_{\mathbf{y}}$ and its gradient through f , i.e., $\hat{\mathbf{y}}_{\mathbf{x}}$ have their roles switched if applied to the conjugate function f^* .

2.1.8 A practical implication

A practical implication of the above is when you need to find stationary point for an dual function $f^*(\cdot)$:

$$\hat{\mathbf{y}}_{\mathbf{x}} = \arg \max_{\mathbf{y}} [\mathbf{x}^\top \mathbf{y} - f^*(\mathbf{y})] \quad (50)$$

and you know by definition, the following is true:

$$\begin{aligned} \nabla_{\mathbf{y}} f^*(\hat{\mathbf{y}}_{\mathbf{x}}) - \mathbf{x} &= 0 \\ \nabla_{\mathbf{y}} f^*(\hat{\mathbf{y}}_{\mathbf{x}}) &= \mathbf{x} \\ &= \hat{\mathbf{x}}_{\mathbf{y}} \end{aligned} \quad (51)$$

one way to solve for $\hat{\mathbf{y}}_{\mathbf{x}}$ is by taking the inverse of the Jacobian matrix:

$$\hat{\mathbf{y}}_{\mathbf{x}} = (\nabla_{\mathbf{y}} f^*)^{-1}(\hat{\mathbf{x}}_{\mathbf{y}}) \quad (52)$$

However, using the property in Eq.(47), we can compute $\hat{\mathbf{y}}_{\mathbf{x}}$ simply by:

$$\hat{\mathbf{y}}_{\mathbf{x}} = \nabla_{\mathbf{x}} f(\hat{\mathbf{x}}_{\mathbf{y}}) \quad (53)$$

2.1.9 derivatives of conjugate function when non-differentiable

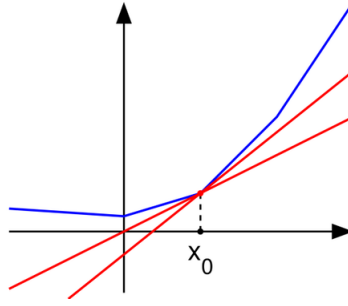
what if $f(\mathbf{x})$ or $f^*(\mathbf{y})$ is **not** differentiable everywhere then the last two lines become:

$$\implies \mathbf{y} \in \partial f(\mathbf{x}) \iff \mathbf{x} \in \partial f^*(\mathbf{y}) \quad (54)$$

in there $\partial f(\mathbf{x})$ is called **sub-differential**

2.1.10 what is sub-differential?

this is used when $f(\mathbf{x})$ is convex, but $f(\mathbf{x})$ is not differentiable at some \mathbf{x}



$$\partial f(x_0) := \{m \in \mathbb{R} \mid f(x) \geq f(x_0) + m(x - x_0) \forall x \in \mathbb{R}\} \quad (55)$$

in words, it means for convex, but non-differentiable function, sub-gradients at point x_0 are the “collection of gradients” where the lines they represent touches the function at $f(x_0)$, but is $\leq f(x)$ everywhere else

Contrary to gradient, i.e., $\frac{df(x)}{dx}$, sub-differential returns a set of gradient values.

If f is differentiable at x , then it can be seen as a special case, $\partial f(x)$ just contains a singleton, $\{\frac{df(x)}{dx}\}$:

$$\partial f(x) := \{m \in \mathbb{R} \mid f(x) \geq f(x_0) + m(x - x_0) \forall x \in \mathbb{R}\} = \left\{ \frac{df(x)}{dx} \right\} \quad (56)$$

2.2 proof for sub differential inverse

In addition to what was described in Section 2, let's look at the inverse from the sub differential perspective. Now we simplify our notation $\hat{\mathbf{y}}_{\mathbf{x}} \rightarrow \hat{\mathbf{y}}$ and $\hat{\mathbf{x}}_{\mathbf{y}} \rightarrow \hat{\mathbf{x}}$ instead:

$$f^*(\hat{\mathbf{y}}) \equiv \sup_{\mathbf{x}} \{\hat{\mathbf{y}}^\top \mathbf{x} - f(\mathbf{x})\} \quad (57)$$

for a particular $\hat{\mathbf{y}}$, assume we find $\hat{\mathbf{x}}$ to be optimum, for which iff:

$$\begin{aligned} \mathbf{0} &\in \hat{\mathbf{y}} - \partial f(\hat{\mathbf{x}}) \\ \implies \hat{\mathbf{y}} &\in \partial f(\hat{\mathbf{x}}) \end{aligned} \quad (58)$$

by substitution $\mathbf{x} = \hat{\mathbf{x}}$ into Eq.(57), (we did not do anything else) we have:

$$f^*(\hat{\mathbf{y}}) \equiv \hat{\mathbf{y}}^\top \hat{\mathbf{x}} - f(\hat{\mathbf{x}}) \quad (59)$$

looking at the equation without picking $\mathbf{y} \equiv \hat{\mathbf{y}}$ (we will add it back into it). For a generic \mathbf{y} :

$$\begin{aligned} f^*(\mathbf{y}) &= \sup_{\mathbf{x}} \{\mathbf{y}^\top \mathbf{x} - f(\mathbf{x})\} \\ &\geq \mathbf{y}^\top \hat{\mathbf{x}} - f(\hat{\mathbf{x}}) \quad \hat{\mathbf{x}} \text{ optimized for } \hat{\mathbf{y}}, \text{ not for generic } \mathbf{y} \\ &= (\mathbf{y} - \hat{\mathbf{y}})^\top \hat{\mathbf{x}} - f(\hat{\mathbf{x}}) + \hat{\mathbf{y}}^\top \hat{\mathbf{x}} \quad \text{add and subtract } \hat{\mathbf{y}}^\top \hat{\mathbf{x}} \\ &= f^*(\hat{\mathbf{y}}) + (\mathbf{y} - \hat{\mathbf{y}})^\top \hat{\mathbf{x}} \quad \text{from Eq.(59)} \quad f^*(\hat{\mathbf{y}}) \equiv \hat{\mathbf{y}}^\top \hat{\mathbf{x}} - f(\hat{\mathbf{x}}) \\ &= f^*(\hat{\mathbf{y}}) + \langle \hat{\mathbf{x}}, \mathbf{y} - \hat{\mathbf{y}} \rangle \end{aligned} \quad (60)$$

then, the moral of the story is that by looking at:

$$f^*(\mathbf{y}) \geq f^*(\hat{\mathbf{y}}) + \langle \hat{\mathbf{x}}, \mathbf{y} - \hat{\mathbf{y}} \rangle \quad (61)$$

since we know $f^*(\mathbf{y})$ must be convex and using the definition of convex function, we know that it must be the case where $\hat{\mathbf{x}} \in \partial_{\mathbf{y}} f^*(\hat{\mathbf{y}})$
we have shown that:

$$\hat{\mathbf{y}} \in \partial f(\hat{\mathbf{x}}) \implies \hat{\mathbf{x}} \in \partial f^*(\hat{\mathbf{y}}) \quad (62)$$

2.3 example of Fenchel/conjugate for $f(\cdot) \equiv \|\cdot\|$

2.3.1 $f \equiv \|\cdot\|$ is a vector norm

$$\begin{aligned} f(\mathbf{x}) &= \|\mathbf{x}\| \\ f^*(\mathbf{y}) &= \sup_{\mathbf{x}} \{\mathbf{y}^\top \mathbf{x} - \|\mathbf{x}\|\} \\ &\equiv \sup_{\mathbf{x}} \{\langle \mathbf{y}, \mathbf{x} \rangle - \|\mathbf{x}\|\} \end{aligned} \quad (63)$$

what is $f^*(\mathbf{y})$? It actually depends on the sign of $\langle \mathbf{y}, \mathbf{x} \rangle - \|\mathbf{x}\|$, and interestingly, we can obtain this through dual norm $\|\mathbf{y}\|_*$ (for its computation, $\|\mathbf{x}\| \leq 1$):

1. case $\|\mathbf{y}\|_* > 1$

using the definition of dual norm:

$$\begin{aligned} \|\mathbf{y}\|_* &= \sup_{\|\mathbf{x}\| \leq 1} \langle \mathbf{y}, \mathbf{x} \rangle \\ \|\mathbf{y}\|_* > 1 &\implies \sup_{\|\mathbf{x}\| \leq 1} \langle \mathbf{y}, \mathbf{x} \rangle > 1 \\ &\implies \exists \mathbf{x} \text{ s.t } \|\mathbf{x}\| \leq 1 : \underbrace{\langle \mathbf{y}, \mathbf{x} \rangle}_{>1} > \underbrace{\|\mathbf{x}\|}_{\leq 1} \end{aligned} \quad (64)$$

therefore, when we apply this \mathbf{y} s.t., $\|\mathbf{y}\|_* > 1$ to unconstrained \mathbf{x} in $\sup_{\mathbf{x}} \{\langle \mathbf{y}, \mathbf{x} \rangle - \|\mathbf{x}\|\}$:

we can rewrite:

$$\mathbf{x} \rightarrow t\mathbf{x} \quad (65)$$

i.e., now \mathbf{x} can have any norm, which is the case of $f^*(\mathbf{y}) \equiv \sup_{\mathbf{x}} \{\langle \mathbf{y}, \mathbf{x} \rangle - \|\mathbf{x}\|\}$, i.e., \mathbf{x} is unbounded:

$$\langle \mathbf{y}, t\mathbf{x} \rangle - \|t\mathbf{x}\| = t(\langle \mathbf{y}, \mathbf{x} \rangle - \|\mathbf{x}\|) \quad (66)$$

$t \rightarrow \infty \implies \|\mathbf{y}\|_*$ to be unbounded, so we have shown half of the dual norm:

$$f^*(\mathbf{y}) = \mathbb{1}_{\|\mathbf{x}\|_* \leq 1}(\mathbf{x}) = \begin{cases} ? & \|\mathbf{y}\|_* \leq 1 \\ +\infty & \|\mathbf{y}\|_* > 1 \end{cases}. \quad (67)$$

2. case $\|\mathbf{y}\|_* \leq 1$

using the definition of dual norm:

$$\begin{aligned} \|\mathbf{y}\|_* &= \sup_{\|\mathbf{x}\| \leq 1} \langle \mathbf{y}, \mathbf{x} \rangle \\ \|\mathbf{y}\|_* \leq 1 &\implies \sup_{\|\mathbf{x}\| \leq 1} \langle \mathbf{y}, \mathbf{x} \rangle \leq 1 \\ &\implies \forall \mathbf{x} \text{ s.t. } \|\mathbf{x}\| \leq 1 : \langle \mathbf{y}, \mathbf{x} \rangle - \|\mathbf{x}\| \leq 0 \end{aligned} \quad (68)$$

since the term inside $\sup_{\mathbf{x}} \{\langle \mathbf{y}, \mathbf{x} \rangle - \|\mathbf{x}\|\}$ is ≤ 0 , and therefore $\sup_{\mathbf{x}}$ makes it 0.

Combined the two, we have:

$$f^*(\mathbf{y}) = \mathbb{1}_{\|\mathbf{y}\|_* \leq 1}(\mathbf{y}) = \begin{cases} 0 & \|\mathbf{y}\|_* \leq 1 \\ +\infty & \|\mathbf{y}\|_* > 1 \end{cases}. \quad (69)$$

In general, indicator function of a convex set:

$$\mathbb{1}_{\mathcal{X}}(\mathbf{x}) = \begin{cases} 0, & \mathbf{x} \in \mathcal{X} \\ +\infty & \text{otherwise} \end{cases} \quad (70)$$

Also, it's important to note that $f(\mathbf{x}) = \|\mathbf{x}\|$ is un-constrained, but $f^*(\mathbf{y})$ is constrained to \mathcal{Y} , i.e., a unit ball of a norm $\|\cdot\| = \{\mathbf{y} : \|\mathbf{y}\| \leq 1\}$.

2.3.2 $f = \|\cdot\|^2$ vector norm

$$f(\mathbf{x}) = \frac{1}{2}\|\mathbf{x}\|^2 \implies f^*(\mathbf{y}) = \frac{1}{2}\|\mathbf{y}\|_*^2 \quad (71)$$

1. show $f^*(\mathbf{y}) \leq \frac{1}{2}\|\mathbf{y}\|_*^2$:

$$\begin{aligned} \mathbf{y}^\top \mathbf{x} &\leq \|\mathbf{y}\|_* \|\mathbf{x}\| \quad \forall \mathbf{x} \\ \implies \mathbf{y}^\top \mathbf{x} - \frac{1}{2}\|\mathbf{x}\|^2 &\leq \|\mathbf{y}\|_* \|\mathbf{x}\| - \frac{1}{2}\|\mathbf{x}\|^2 \quad \text{both sides } - \frac{1}{2}\|\mathbf{x}\|^2 \end{aligned} \quad (72)$$

R.H.S is quadratic function in terms of $\|\mathbf{x}\|$. After solving, it has max occur at $\|\mathbf{x}\| = \|\mathbf{y}\|_*$, substituting into:

$$\begin{aligned} \max_{\|\mathbf{x}\|} \left(\|\mathbf{y}\|_* \|\mathbf{x}\| - \frac{1}{2}\|\mathbf{x}\|^2 \right) &= \|\mathbf{y}\|_*^2 - \frac{1}{2}\|\mathbf{y}\|_*^2 \\ &= \frac{1}{2}\|\mathbf{y}\|_*^2 \\ \implies f^*(\mathbf{y}) = \mathbf{y}^\top \mathbf{x} - \frac{1}{2}\|\mathbf{x}\|^2 &\leq \frac{1}{2}\|\mathbf{y}\|_*^2 \quad \forall \mathbf{x} \end{aligned} \quad (73)$$

2. $f^*(\mathbf{y}) \geq \frac{1}{2}\|\mathbf{y}\|_*^2$

Let \mathbf{x} be any chosen vector with $\mathbf{y}^\top \mathbf{x} = \|\mathbf{y}\|_* \|\mathbf{x}\|$ (think this as when norm on the R.H.S is fixed, then one may choose the value of \mathbf{y} and \mathbf{x} to change the “direction”. In $\|\cdot\|_2$ case, they should have $\cos(\theta) = 0$). W.l.o.g, we then scale so that $\|\mathbf{x}\| = \|\mathbf{y}\|_*$. then we have, for this \mathbf{x} :

$$\begin{aligned} \|\mathbf{x}\| &= \|\mathbf{y}\|_* \\ \implies \mathbf{y}^\top \mathbf{x} &= \|\mathbf{x}\|^2 \quad \because \mathbf{y}^\top \mathbf{x} = \|\mathbf{y}\|_* \|\mathbf{x}\| \\ &= \frac{1}{2}\|\mathbf{x}\|^2 + \frac{1}{2}\|\mathbf{y}\|_*^2 \\ \implies \mathbf{y}^\top \mathbf{x} - \frac{1}{2}\|\mathbf{x}\|^2 &= \frac{1}{2}\|\mathbf{y}\|_*^2 \\ \implies f^*(\mathbf{y}) &\geq \frac{1}{2}\|\mathbf{y}\|_*^2 \end{aligned} \quad (74)$$

References

- [1] Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro, “The implicit bias of gradient descent on separable data,” *The Journal of Machine Learning Research*, vol. 19, no. 1, pp. 2822–2878, 2018.