

Plan { Optimality Condition on a manifold  
 Convergence of Riemannian GD

---

We discuss optimality using an example:

$$\textcircled{1} \begin{cases} \min_{X \in \mathbb{R}^{n \times p}} f(x) = \text{tr}(X^T A X) & A^T = A \in \mathbb{R}^{n \times n} \\ \text{s.t. } X^T X = I \end{cases}$$

$$\Leftrightarrow \textcircled{2} \min_{X \in M} f(x), \quad M = \text{St}(n, p) = \{X \in \mathbb{R}^{n \times p} : X^T X = I\}$$

I. Necessary but not sufficient optimality condition for  $\textcircled{1}$   
 can be derived from saddle point of Lagrangian

$$L(x, \Lambda) = f(x) - \langle \Lambda, X^T X - I \rangle$$

$$X \in \mathbb{R}^{n \times p} \quad \Lambda \in \mathbb{R}^{p \times p} \quad \langle U, V \rangle = \sum_i \sum_j U_{ij} V_{ij}$$

Regard  $f(x)$  as a function on  $\mathbb{R}^{n \times p}$

$$\begin{aligned} \frac{\partial L}{\partial x} &= \frac{\partial}{\partial x} f(x) - \frac{\partial}{\partial x} \langle \Lambda, X^T X - I \rangle \\ &= \frac{\partial}{\partial x} \langle x, Ax \rangle - \frac{\partial}{\partial x} \langle \Lambda, X^T X \rangle \\ &= 2Ax - (X\Lambda + X\Lambda^T) \end{aligned}$$

How to calculate  $\frac{\partial}{\partial x} \langle \Lambda, X^T X \rangle$ :  $\text{tr}(A^T B) = \langle A, B \rangle$

$$1) \langle \Lambda, Y^T X \rangle = \text{tr}(\Lambda^T Y^T X) = \text{tr}((Y\Lambda)^T X) = \langle Y\Lambda, X \rangle$$

$$\Rightarrow \frac{\partial}{\partial x} \langle \Lambda, Y^T X \rangle = Y\Lambda$$

$$2) \langle \Lambda, Y^T X \rangle = \text{tr}(X^T Y \Lambda) = \text{tr}(\Lambda X^T Y) = \langle X \Lambda^T, Y \rangle$$

$\downarrow$   
 $\text{tr}(ABC) = \text{tr}(CAB)$

$$\Rightarrow \frac{\partial}{\partial Y} \langle \Lambda, Y^T X \rangle = X \Lambda^T$$

$$3) \frac{\partial}{\partial X} \langle \Lambda, X^T X \rangle = \frac{\partial}{\partial X} \langle \Lambda, Y^T X \rangle \Big|_{Y=X} + \frac{\partial}{\partial Y} \langle \Lambda, Y^T X \rangle \Big|_{Y=X}$$

$$= X \Lambda + X \Lambda^T$$

$$L(x, \Lambda) = f(x) - \langle \Lambda, X^T X - I \rangle$$

Saddle Point

$$\begin{cases} \frac{\partial L}{\partial X} = \frac{\partial}{\partial X} f(x) - (X \Lambda + X \Lambda^T) = 0 \\ \frac{\partial L}{\partial \Lambda} = X^T X - I = 0 \end{cases}$$

II. Necessary but not sufficient optimality condition for ②

- a manifold  $M \subset \mathcal{E}$   
 $\begin{matrix} \parallel & \parallel \\ \text{St}(n, p) & \mathbb{R}^{n \times p} \end{matrix}$

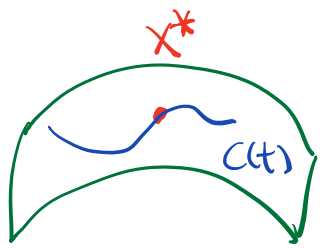
- a curve  $C: (-\epsilon, \epsilon) \rightarrow M$   
 $t \mapsto C(t)$

To understand/calculate  $C'(t)$ , we identify it with a natural extension

$$C: \mathbb{R} \rightarrow \mathcal{E}$$

$$t \mapsto C(t)$$

- If  $x_* \in M$  is the minimizer of  $f(x)$  on  $M$ , then consider any curve  $c(t)$  on  $M$  with  $c(0) = x_*$ .



$f \circ c(t)$  has a minimizer at  $t=0$

$$\Rightarrow \frac{d}{dt}[f \circ c] \Big|_{t=0} = 0$$

- Recall the definition of **differential** of  $f: M \rightarrow \mathbb{R}$  is a linear map

$$Df(x): T_x M \rightarrow \mathbb{R}$$

$$v \mapsto Df(x)[v] \stackrel{\text{def}}{=} \frac{d}{dt} f(\gamma(t)) \Big|_{t=0}$$

where  $\gamma(t)$  is any curve on  $M$  with  $\begin{cases} \gamma(0) = x \\ \gamma'(0) = v \end{cases}$

- Theorem (Necessary Optimality Condition on  $M$ )

$$(a) \begin{cases} \frac{d}{dt}[f \circ c] \Big|_{t=0} = 0 \\ \text{for any curve } c(t) \text{ on } M \text{ with } c'(0) = x_* \end{cases}$$

$$\Leftrightarrow (b) Df(x_*)[v] = 0, \forall v \in T_{x_*} M$$

$$\Leftrightarrow (c) \text{grad} f(x_*) = 0$$

Proof: (b)  $\Leftrightarrow$  (c) def of Riemannian Grad

$$g_{x_*}(\text{grad} f(x_*), v) = Df(x_*)[v], \forall v \in T_{x_*} M$$

Since  $\text{grad} f(x_*) \in T_{x_*} M$ , we pick  $v = \text{grad} f(x_*)$ ,

$$\text{then } g_{x_*}(\text{grad} f(x_*), \text{grad} f(x_*)) = 0$$

$$\Rightarrow \text{grad} f(x_*) = 0$$

positive-definiteness of metric

$$(a) \Leftrightarrow (b) \quad (a) \begin{cases} \frac{d}{dt}[f \circ c] \Big|_{t=0} = 0 \\ \text{for any curve } c(t) \text{ on } M \text{ with } c(0) = x_* \end{cases}$$

The general definition of tangent space of a manifold  $M \subseteq \mathbb{E}$

$$T_x M = \{c'(0) \mid c: (-\epsilon, \epsilon) \rightarrow M \text{ is smooth \& } c(0) = x\}$$

For any curve  $\gamma(t)$  with  $\gamma(0) = x_*$ , let  $v = \gamma'(0)$

$$\text{then } Df(x)[v] \stackrel{\text{def}}{=} \frac{d}{dt} f(\gamma(t)) \Big|_{t=0} = 0$$

arbitrariness of such a curve  $\Leftrightarrow \forall v \in T_{x_*} M$

• So the optimality on  $M$  is  $\text{grad} f(x_*) = 0$

$$\begin{aligned} g_{x_*}(\text{grad} f(x_*), v) &= Df(x_*)[v] = D\bar{f}(x_*)[v] \\ &= \langle \nabla \bar{f}(x_*), v \rangle \end{aligned}$$

$$\begin{array}{ll} f: M \rightarrow \mathbb{R} & \text{has an extension } \bar{f}: \mathbb{R}^{n \times p} \rightarrow \mathbb{R} \\ x \mapsto \text{tr}(x^T A x) & x \mapsto \text{tr}(x^T A x) \end{array}$$

In particular, if we choose  $g_x(\cdot)$  to be  $\langle \cdot, \cdot \rangle$   
 then  $\text{grad} f(x_*)$  is Proj of  $\nabla \bar{f}(x_*)$   
 on  $T_{x_*}M$

Previously, we have computed

$$\forall Y \in \mathbb{R}^{n \times p}, P_{T_x M}(Y) = (I - xx^T)Y + \text{skew}(x^T Y)$$

$$T_x \text{St}(n, p) = \{ V \in \mathbb{R}^{n \times p} : x^T V + V^T x = 0 \} \subseteq \mathbb{R}^{n \times p}$$

$$\begin{aligned} \text{So } \text{grad} f(x) &= (I - xx^T) \nabla \bar{f}(x) + \frac{x^T \nabla \bar{f}(x) - \nabla \bar{f}(x)^T x}{2} \\ &= (I - xx^T) 2Ax + x^T Ax - x^T Ax \\ &= (I - xx^T) 2Ax \quad \textcircled{2} \end{aligned}$$

Let's compare it to

$$\text{Optimality via } L(x, \Lambda) = f(x) - \langle \Lambda, x^T x - I \rangle$$

$$\textcircled{1} \quad \text{Saddle Point} \begin{cases} \frac{\partial L}{\partial x} = \frac{\partial}{\partial x} f(x) - (x\Lambda + x\Lambda^T) = 0 \\ \frac{\partial L}{\partial \Lambda} = x^T x - I = 0 \end{cases}$$

$$\text{(Claim : } \textcircled{1} \Leftrightarrow \textcircled{2} \text{ (grad } f(x) = P_{T_x M}(\nabla \bar{f}(x)) = 0 \text{))}$$

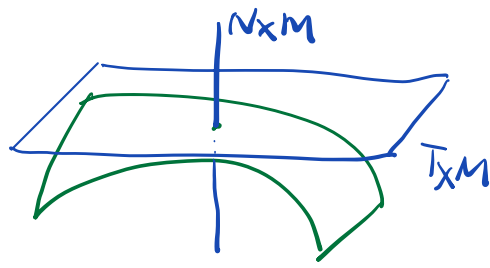
Proof: Previously we derived the normal space

$$N_x \text{St}(n, p) = \{ X S : S^T = S, S \in \mathbb{R}^{p \times p} \}$$

$T_x M$  is a subspace in  $\mathcal{E}$

$N_x M$  is the orthogonal complement of  $T_x M$

and  $\mathcal{E} = T_x M \oplus N_x M$



$$\frac{\partial}{\partial x} f(x) - (x\Lambda + x\Lambda^T) = 0$$

$$\Leftrightarrow \frac{\partial}{\partial x} f(x) = X \underbrace{(\Lambda + \Lambda^T)}_S \in N_x \text{St}(n, p)$$

$$\Leftrightarrow \text{Proj of } \frac{\partial}{\partial x} f(x) \text{ onto } T_x M \text{ is } 0$$

$$\Leftrightarrow \textcircled{2} \quad \text{because } \frac{\partial}{\partial x} f(x) \text{ in } \textcircled{1} \text{ is } \nabla \bar{f}(x) \text{ in } \textcircled{2}$$

• Remark & Observation :

1) When using Euclidean metric  $g_x(\cdot, \cdot) = \langle \cdot, \cdot \rangle$

$$\text{grad} f(x) = P_{T_x M} \left( \frac{\partial}{\partial x} f(x) \right), \text{ so } \textcircled{1} \Leftrightarrow \textcircled{2}$$

2) If using a generic metric  $g$ , then

$$\text{grad} f(x) \neq P_{T_x M} \left( \frac{\partial}{\partial x} f(x) \right).$$

3) Intuitively, optimality conditions  $\textcircled{1}$  &  $\textcircled{2}$

should be equivalent, regardless of metric.

But if  $\text{grad} f(x) \neq P_{T_x M} \left( \frac{\partial}{\partial x} f(x) \right)$ ,

why  $\textcircled{1} \Leftrightarrow \textcircled{2}$ ?

Answer:  $\textcircled{1} \Leftrightarrow P_{T_x M} \left( \frac{\partial}{\partial x} f(x) \right) = 0 \Leftrightarrow \frac{\partial}{\partial x} f(x) \in N_x M$

$$\textcircled{2} \Leftrightarrow \text{grad} f(x) = 0$$

$$\textcircled{1} \Rightarrow g_x(\text{grad} f(x), v) = \langle \nabla \bar{f}(x), v \rangle = 0, \forall v \in T_x M$$

$$\Rightarrow g_x(\text{grad} f(x), v) = 0$$

$$\Rightarrow g_x(\text{grad} f(x), \text{grad} f(x)) = 0 \Rightarrow \text{grad} f(x) = 0$$

$$\textcircled{2} \Rightarrow \forall v \in T_x M, g_x(\text{grad} f(x), v) = 0$$

$$\Rightarrow \langle \nabla \bar{f}(x), v \rangle = g_x(\text{grad} f(x), v) = 0$$

$$\Rightarrow \nabla \bar{f}(x) \in N_x M$$

$$\Rightarrow \text{Proj of } \nabla \bar{f}(x) \text{ onto } T_x M \text{ is } 0$$

---

Riemannian Gradient Descent for  $\min_{x \in M} f(x)$

$$x_{k+1} = R_{x_k}(-\eta_k \text{grad} f(x_k))$$

We expect  $x_k$  converge to  $x_*$  s.t.  $\text{grad} f(x_*) = 0$

↓  
just a critical point  
not necessarily a minimizer

Step Size rule/method:

1) a constant step size  $\eta_k = \eta$

2) "optimal" step size

$$\eta_k = \underset{t}{\text{argmin}} h(t) = R_{x_k}(-t \text{grad} f(x_k))$$

### 3) Line search by backtracking

Start with step size  $t_0$ , iteratively reduce it to  $t_i = \rho t_{i-1}$  for some  $\rho \in (0, 1)$  until some conditions are satisfied.

Ideas/Steps to show the convergence:

① Assume  $f(x) \geq D, \forall x \in M$

② Show sufficient decrease

$$f(x_k) - f(x_{k+1}) \geq c \|\text{grad} f(x_k)\|^2$$

only possible with  $\left\{ \begin{array}{l} \text{certain step sizes} \\ \text{assumptions and } f \& (M, g) \end{array} \right.$

③ Then we can show  $\lim_{k \rightarrow \infty} \|\text{grad} f(x_k)\| = 0$

$$0 \leq \sum_{k=0}^N [f(x_k) - f(x_{k+1})] = f(x_0) - f(x_N) \leq f(x_0) - D$$

$\Rightarrow \sum_{k=0}^{\infty} [f(x_k) - f(x_{k+1})]$  is finite

$\Rightarrow c \|\text{grad} f(x_k)\|^2 \leq f(x_k) - f(x_{k+1}) \rightarrow 0, k \rightarrow \infty$



$$\begin{aligned}
f(x_0) - D &\geq f(x_0) - f(x_N) \\
&= \sum_{k=0}^{N-1} [f(x_k) - f(x_{k+1})] \\
&\geq N c \min_{0 \leq k \leq N-1} \|\text{grad} f(x_k)\|^2
\end{aligned}$$

$$\Rightarrow \min_{0 \leq k \leq N-1} \|\text{grad} f(x_k)\| \leq \sqrt{\frac{f(x_0) - D}{c}} \frac{1}{\sqrt{N}}$$

④ How to get sufficient decrease

$$f(x_k) - f(x_{k+1}) \geq c \|\text{grad} f(x_k)\|^2$$

1) For special  $f$ ,  $R_x$ ,

$$v \in T_x M, f(R_x(v)) \leq f(x) + \langle \text{grad} f(x), v \rangle + \frac{L}{2} \|v\|^2$$

2) With constant step size  $\eta_k = \frac{1}{L}$ , can show

$$f(x_k) - f(x_{k+1}) \geq \frac{1}{2L} \|\text{grad} f(x_k)\|^2$$

Proof:  $v = -\frac{1}{L} \text{grad} f(x_k)$