

535520: Optimization Algorithms

Lecture 11 – Mirror Descent and Newton's Method

Ping-Chun Hsieh (謝秉均)

November 25, 2024

Announcement

- ▶ Lecture today (11/25): 12:20pm-2:20pm
- ▶ 10min extension for the next three lectures

This Lecture

1. Mirror Descent

2. Newton's Method

- Reading Material:
 - Amir Beck and Marc Teboulle, “Mirror descent and nonlinear projected subgradient methods for convex optimization,” Operations Research Letters, 2003.
 - Chapters 3 and 6 of Nocedal & Wright’s textbook “Numerical Optimization”
 - Chapter 9 of Stephen Boyd’s textbook “Convex Optimization”

Scipy.Optimize in Python?

Optimization and root finding (`scipy.optimize`)

SciPy `optimize` provides functions for minimizing (or maximizing) objective functions, possibly subject to constraints. It includes solvers for nonlinear problems (with support for both local and global optimization algorithms), linear programming, constrained and nonlinear least-squares, root finding, and curve fitting.

Common functions and objects, shared across different solvers, are:

`show_options([solver, method, disp])` Show documentation for additional options of optimization solvers.

`OptimizeResult` Represents the optimization result.

`OptimizeWarning`

Local (multivariate) optimization

`minimize(fun, x0[, args, method, jac, hess, ...])`

Minimization of scalar function of one or more variables.

The `minimize` function supports the following methods:

`minimize(method='Nelder-Mead')`

`minimize(method='Powell')`

`minimize(method='CG')`

`minimize(method='BFGS')`

`minimize(method='Newton-CG')`

`minimize(method='L-BFGS-B')`

`minimize(method='TNC')`

`minimize(method='COBYLA')`

Newton's

↓

Quasi-Newton

Review: Mirror Descent (MD)

Key Idea: Generalize the L_2 proximal term to other distance measures!

$$x_{t+1} = \arg \min_{x \in C} \left\{ f(x_t) + \nabla f(x_t)^\top (x - x_t) + \frac{1}{\eta_t} D_\phi(x \| x_t) \right\}$$

first-order approximation Bregman divergence

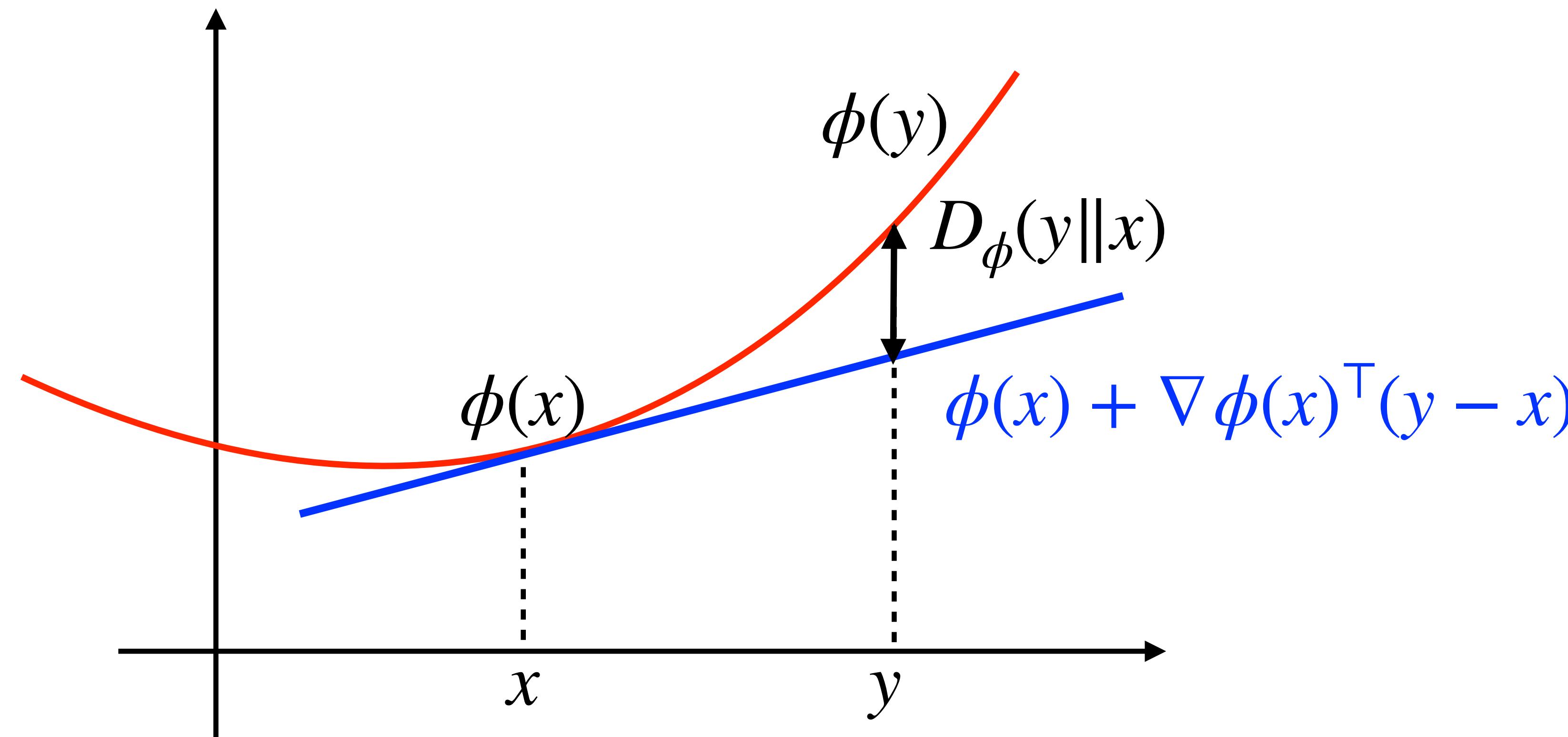
where $D_\phi(y \| x) := \phi(y) - \phi(x) - \nabla \phi(x)^\top (y - x)$

(with respect to $\phi(\cdot)$ strictly convex and differentiable)

-
- ▶ **Remark:** Bregman divergence is meant to capture "local geometry" of objective function
 - ▶ **Remark:** Bregman divergence is NOT symmetric and hence not a metric
 - ▶ How about MD in the unconstrained cases?

Review: Bregman Divergence

$$D_\phi(y\|x) = \underline{\phi(y)} - \left(\underline{\phi(x)} + \underline{\nabla\phi(x)^\top(y-x)} \right), \text{ where } \phi(\cdot) \text{ is strictly convex}$$



Review: Basic Properties of Bregman Divergence

Let $\phi : X \rightarrow \mathbb{R}$ be a **strictly convex** and **differentiable** function

- ▶ 1. **Non-negativity**: $D_\phi(y\|x) \geq 0$
- ▶ 2. **Distance between a point to itself is zero**: $D_\phi(y\|x) = 0$ if and only if $x = y$
- ▶ 3. **Convexity**: $D_\phi(y\|x)$ is convex in y (but not necessarily in x)
- ▶ 4. **Symmetry not guaranteed**: $D_\phi(y\|x) \neq D_\phi(x\|y)$ in general

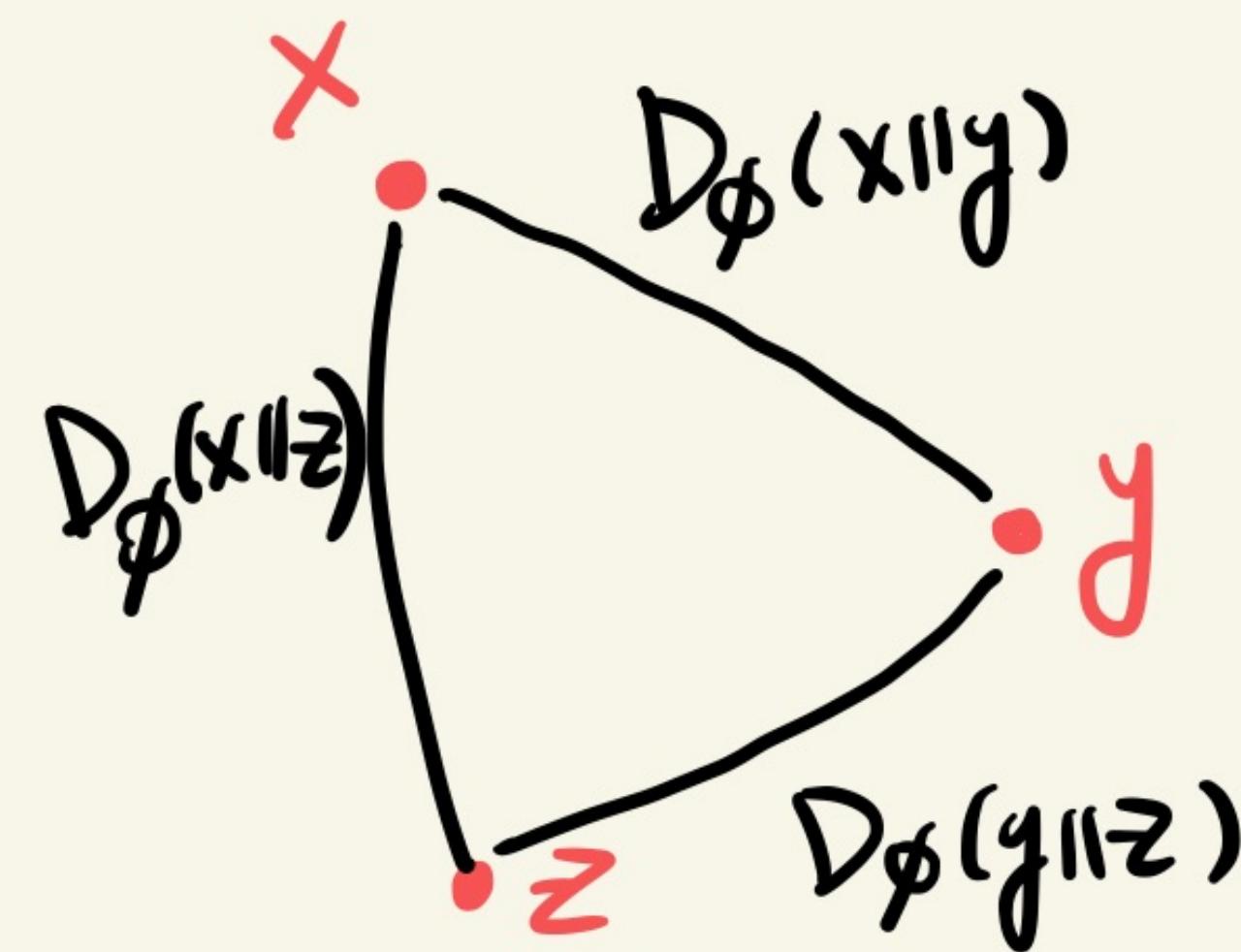
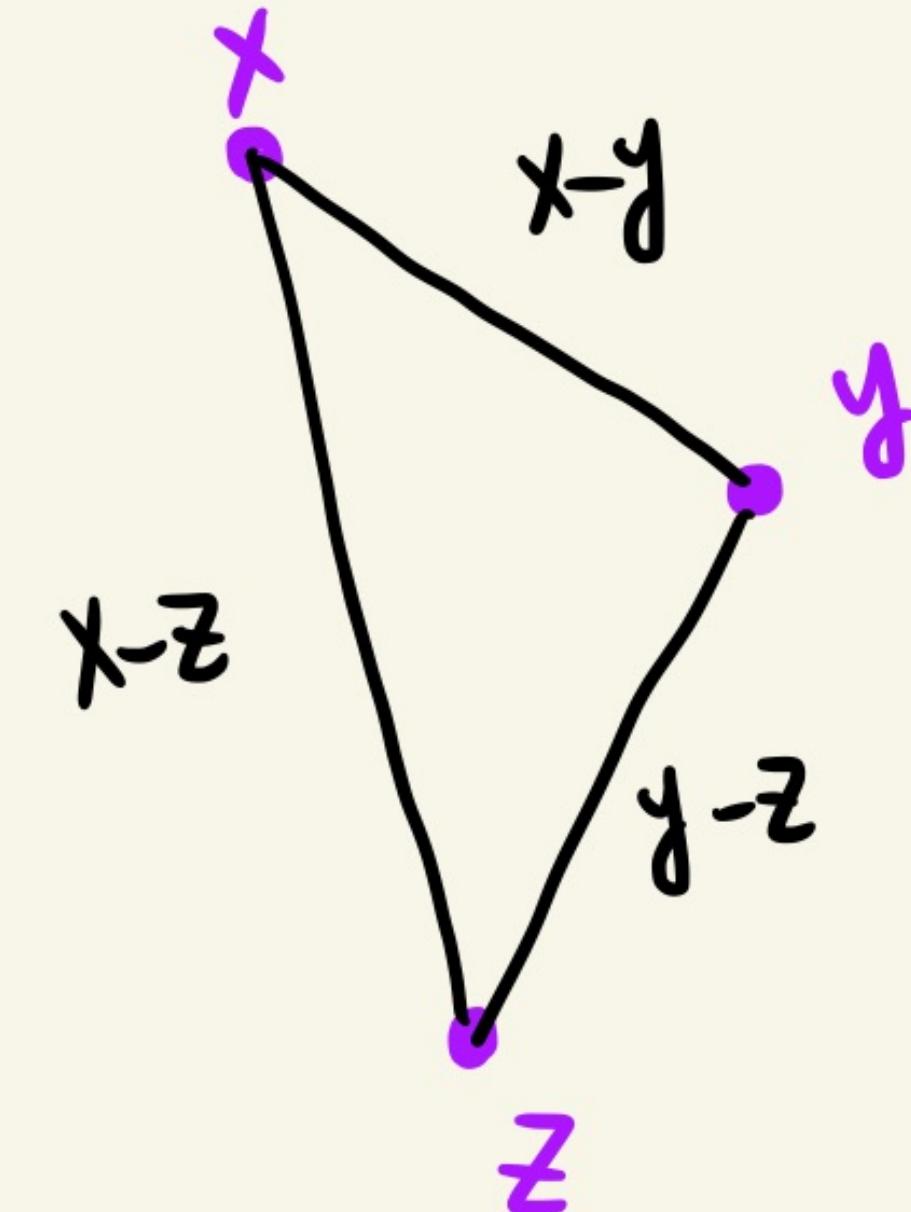
Bregman Divergence: Three-Point Lemma

Euclidean distance

$$\|x-z\|^2 = \|x-y\|^2 + \|y-z\|^2 - 2(y-z)^T(x-y)$$

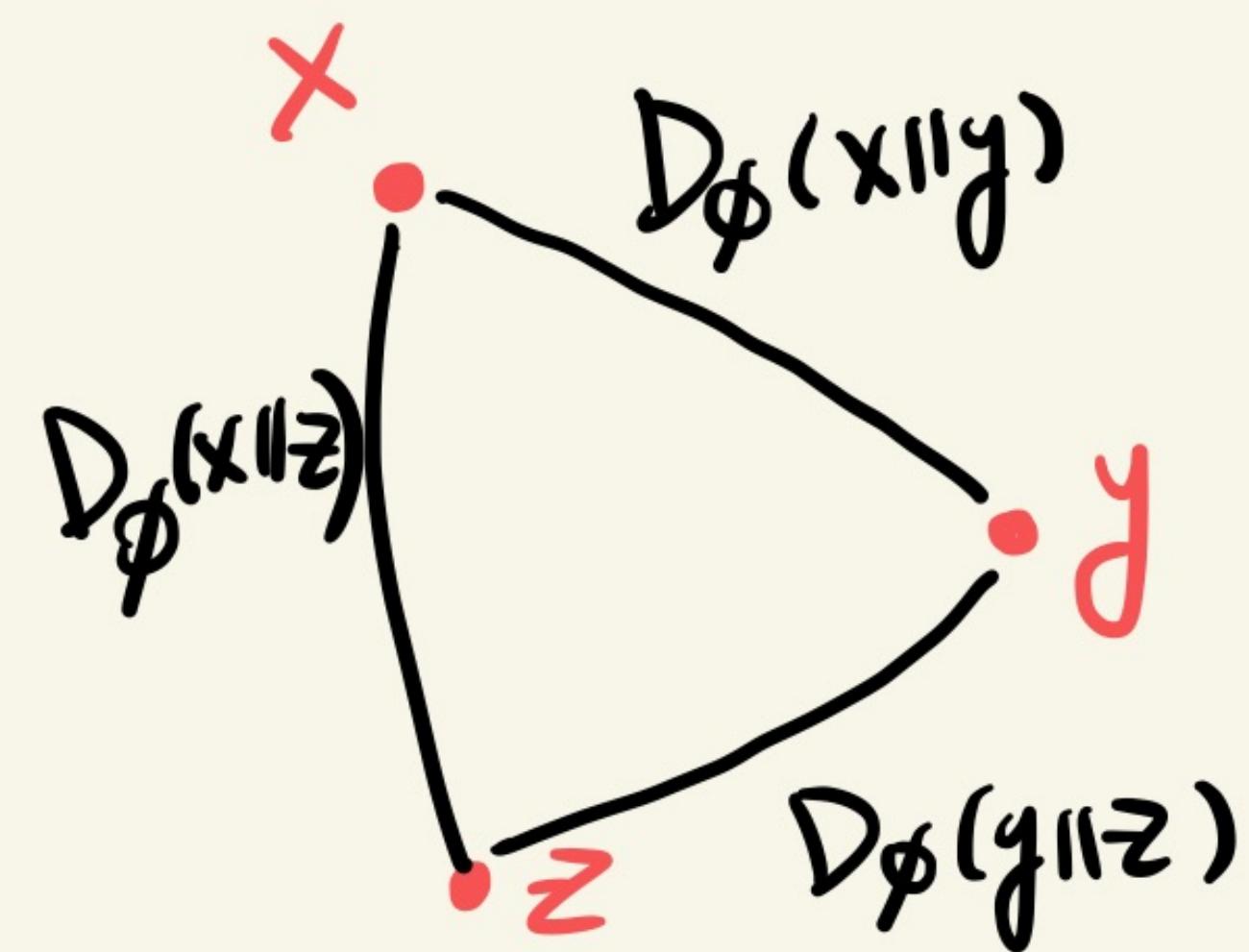
Bregman divergence

$$D_\phi(x||z) = D_\phi(x||y) + D_\phi(y||z) - (\nabla \phi(z) - \nabla \phi(y))^\top (x-y)$$



Proof of Three-Point Lemma

$$\begin{aligned}
 & D_\phi(x||y) + D_\phi(y||z) - D_\phi(x||z) \\
 = & \left(\phi(x) - \phi(y) - \nabla\phi(y)^T(x-y) \right) \\
 + & \left(\phi(y) - \phi(z) - \nabla\phi(z)^T(y-z) \right) \\
 - & \left(\phi(x) - \phi(z) - \nabla\phi(z)^T(x-z) \right) \\
 = & (\nabla\phi(z) - \nabla\phi(y))^T(x-y)
 \end{aligned}$$



Another Viewpoint: MD as *Follow-The-Regularized-Leader* in Online Linear Optimization

A Motivating Example of Online Linear Optimization

- ▶ Example: Expert Problem
- ▶ A learner and n experts
- ▶ At each time t , the learner chooses an action probability vector $a_t \in \Delta_n$ (indicating the distribution of expert-following behavior of the learner)
- ▶ Let $y_t \in \mathbb{R}_+^n$ denote the cost vector of following the experts (known after choosing a_t)
- ▶ Goal: Minimize “total cost” = Minimize “regret” (relative to a fixed comparator $a \in \Delta_n$)

$$R_T := \max_{a \in \Delta_n} \left\{ \sum_{t=1}^T (a_t - a)^\top y_t \right\}$$

$$\begin{bmatrix} 2.5 \\ 5 \\ 10 \end{bmatrix}$$

$$\min \sum_{t=1}^T a_t^\top y_t$$

At each time t :

Based on the history $\{y_1, y_2, \dots, y_{t-1}\}$,
select your $a_t \in \Delta_n$

Interesting Fact: $R_T = O(\sqrt{T \log n})$ under properly-designed algorithms!

Suppose $n=3$

$$a_t = \begin{bmatrix} 0.2 \\ 0.5 \\ 0.3 \end{bmatrix}$$

“adversarial”

MD and “Follow-The-Regularized-Leader” (FTRL)

Mirror Descent

Initially:

$$a_1 = \arg \min_{a \in A} \phi(a)$$

For $t \geq 1$:

$$a_{t+1} = \arg \min_{a \in A} \left\{ a^T y_t + \frac{1}{\eta} D_\phi(a \| a_t) \right\}$$

fine-tuning
based on the new observation

↓
fine-tuning based on
the current action a_t

Follow-The-Regularized-Leader

Initially:

$$a_1 = \arg \min_{a \in A} \phi(a)$$

For $t \geq 1$:

$$a_{t+1} = \arg \min_{a \in A} \left\{ \sum_{s=1}^t a^T y_s + \frac{1}{\eta} \phi(a) \right\}$$

(Choosing the action that appears the best in hindsight under regularization)

Interesting Fact: MD and FTRL are equivalent! (under some mild conditions)

Equivalence Between MD and FTRL

MD

$$a_{t+1} = \underset{a \in A}{\operatorname{argmin}} \left\{ a^T y_t + \frac{1}{2} \cdot D_\phi(a \| a_t) \right\}$$

If a_{t+1} is in the *interior* of A , then

By FONC:

$$\nabla F(a_{t+1}) = 0 \quad (\text{and hence } -\eta \cdot y_t = \underline{\nabla \phi(a_{t+1}) - \nabla \phi(a_t)})$$

Therefore, by taking "telescoping sum", we have

$$\nabla \phi(a_{t+1}) = -\eta \cdot \sum_{s=1}^t y_t$$

$$\Rightarrow a_{t+1} = (\nabla \phi)^{-1} \left(-\eta \cdot \sum_{s=1}^t y_t \right)$$

$$\begin{aligned} F(a) &= a^T y_t + \frac{1}{2} \left(\phi(a) - \phi(a_t) \right) \\ &\quad - \nabla \phi(a_t)^T (a - a_t) \end{aligned}$$

$$\nabla F(a) = y_t + \frac{1}{2} \left(\nabla \phi(a) - \nabla \phi(a_t) \right)$$

Equivalence Between MD and FTRL

FTRL

$$a_{t+1} = \underset{a \in A}{\operatorname{argmin}} \left\{ \sum_{s=1}^t a^T y_s + \frac{1}{2} \phi(a) \right\}$$

$G(a)$
||
0

If a_{t+1} is in the interior of A , then

$$\nabla G(a_{t+1}) = 0$$

(Equivalently ,

$$\nabla \phi(a_{t+1}) = -\eta \cdot \sum_{s=1}^t y_s$$

by FONC

Why Not Using “Follow-The-Leader” (FTL)?

Follow-The-Leader: Choosing the action that appears the best in hindsight

$$a_{t+1} = \arg \min_{a \in A} \left\{ \sum_{s=1}^t a^\top y_s \right\}$$

- ▶ FTL (i.e., without a regularizer) suffers from linear regret!

-
- ▶ Counterexample: *action : [-dimensional*

- ▶ Let action set $A = [-1, 1]$ and initial action $a_1 = 0$

- ▶ Cost vectors: $y_1 = 1/2$ and $y_s = (-1)^{s+1}, \forall s > 1$

Counterexample of "Follow-The-Leader"

$$A = [-1, 1], \quad a_1 = 0$$

Total cost up to time t

$$\sum_{s=1}^t a_s^T y_s = t - 1$$

$$y_1 = \frac{1}{2} \Rightarrow a_2 = \underset{a \in A}{\operatorname{argmin}} \quad a^T y_1 = -1$$

$$y_2 = -1 \Rightarrow a_3 = \underset{a \in A}{\operatorname{argmin}} \quad a^T (y_1 + y_2) = +1$$

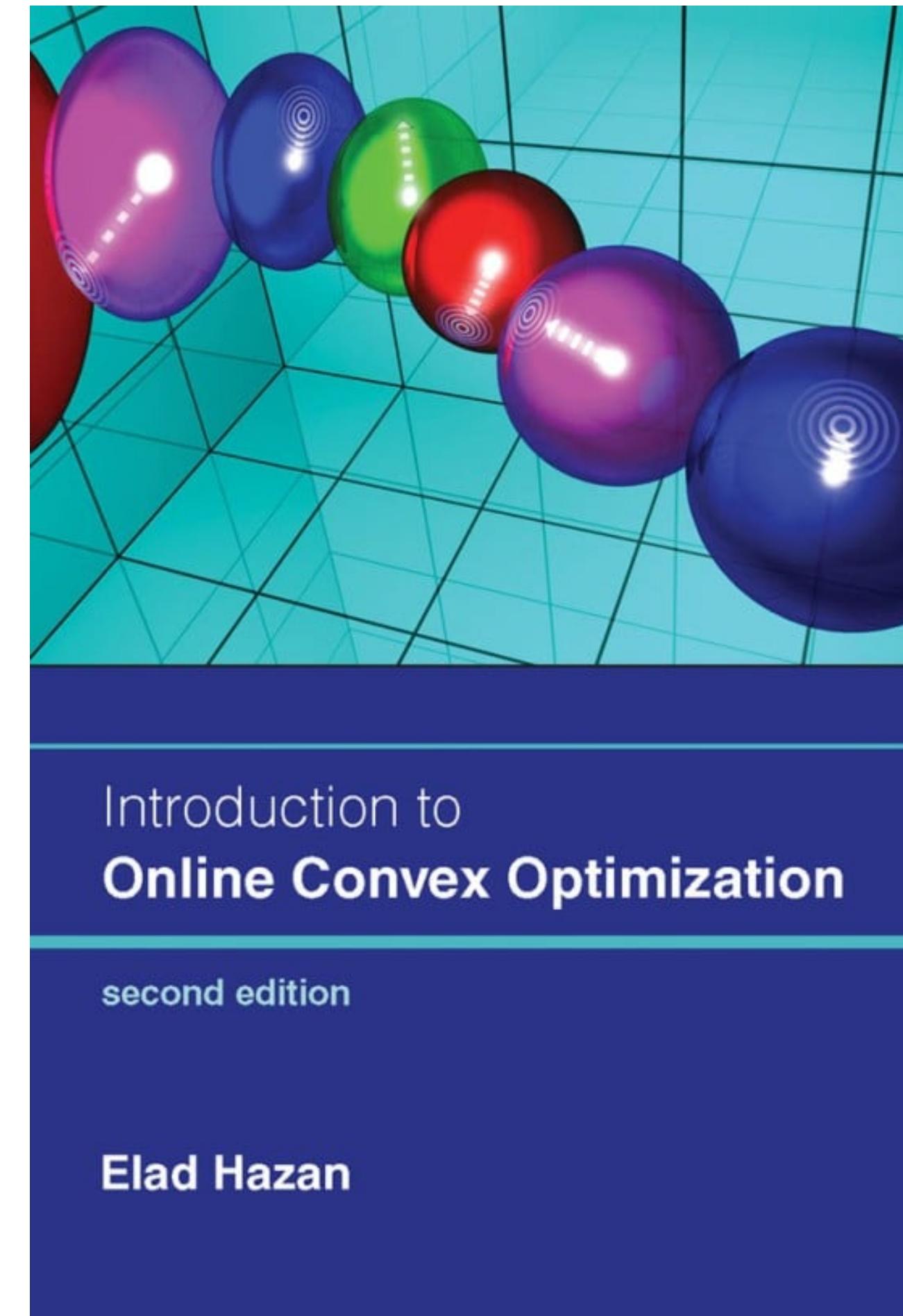
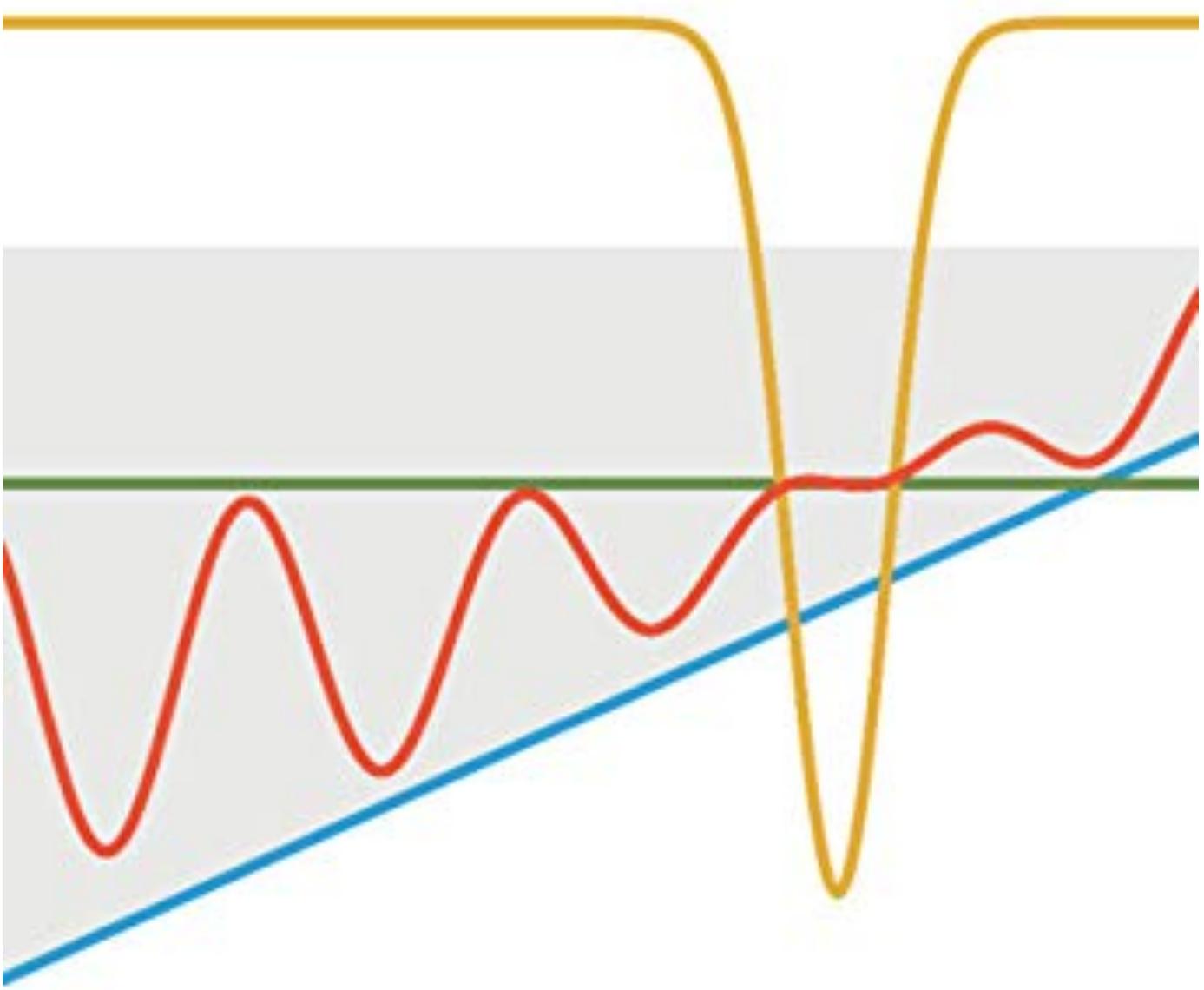
$$y_3 = +1 \Rightarrow a_4 = \underset{a \in A}{\operatorname{argmin}} \quad a^T (y_1 + y_2 + y_3) = -1$$

Hence, $a_t = (-1)^{t+1}$

"Best action in hindsight" is 0 , Regret $R_T = T - 1$
 We know y_1, \dots, y_t

Bandit Algorithms

TOR LATTIMORE
CSABA SZEPESVÁRI



<https://tor-lattimore.com/downloads/book/book.pdf>

<https://arxiv.org/abs/1909.05207>

Convergence of Mirror Descent

Convergence of MD Algorithms

$$\frac{1}{P} + \frac{1}{Q} = 1.$$

Theorem

Suppose the following conditions hold:

① f is convex and L_f -Lipschitz continuous (in the sense that $\|\nabla f(x)\|_F \leq L_f$)

② ϕ is ρ -Strongly convex w.r.t. $\|\cdot\|_P$

Then, MD achieves $f_{best}^t - f^* \leq$

Moreover, by letting

$$R_t = \frac{\sqrt{2\rho \cdot D_{max}}}{L_f} \cdot \frac{1}{\sqrt{t}} = O\left(\frac{1}{\sqrt{t}}\right)$$

$$f_{best}^t - f^* = O\left(\frac{L_f \cdot \sqrt{D_{max}}}{\sqrt{\rho}} \cdot \frac{\log t}{\sqrt{t}}\right)$$

$$\frac{L_f^2}{2\rho} \cdot \sum_{t=0}^T R_t^2$$

step size

$$\sup_{x \in C} D_\phi(x \| x_0) +$$

$$D_{max}$$

Example: Optimization over Probability Simplex (For simplicity, let $X_0 = (\frac{1}{d}, \dots, \frac{1}{d})$)

Suppose we consider: $\min f(x)$, subject to $x \geq 0$, $1^T x = 1$.

Let's compare MD with "Euclidean norm" vs "KL divergence"

① Euclidean norm:

$\phi(x) = \frac{1}{2} \|x\|^2$ is 1-strongly convex w.r.t. $\|\cdot\|_2$

By the definition of Bregman divergence

$$\phi(y) - (\phi(x) + \nabla \phi(x)^T (y-x)) \geq \frac{1}{2} \|y-x\|^2$$

$$D_{\max} := \sup_{x \in C} D_\phi(x \| X_0) = \sup_{x \in C} \frac{1}{2} \|x - (\frac{1}{d}, \dots, \frac{1}{d})\|^2 \leq \frac{1}{2}$$

Hence, $f_{\text{best}}^t - f^* = O\left(\frac{L_f \cdot \sqrt{D_{\max}}}{\sqrt{\rho}} \cdot \frac{\log t}{\sqrt{t}}\right) = O\left(L_{f,2} \cdot \frac{\log t}{\sqrt{t}}\right)$

$$\text{gap}(t+1) \leq \underbrace{\text{const.}}_{\downarrow} \cdot \text{gap}(t) + \cancel{\alpha} \begin{pmatrix} \text{GD} \\ \text{NAG} \end{pmatrix}$$

$$\text{gap}(t) = \frac{1}{t}$$

$$\text{gap}(t) = \beta^t$$

$$\beta^{t+1} \leq (-1) \beta^t + \circ$$

(Cont.)

② KL divergence:

- $\phi(x) = -\sum_{i=1}^d x^{(i)} \log x^{(i)}$ is 1 -strongly convex w.r.t. $\|\cdot\|_1$

$$\begin{aligned}\sup_{x \in C} D_\phi(x \| x_0) &= \sup_{x \in C} \text{KL}(x \| x_0) = \sup_{x \in C} \left(\sum_{i=1}^d x^{(i)} \log x^{(i)} - \sum_{i=1}^d x^{(i)} \cdot \log \frac{1}{n} \right) \\ &= \log n + \sup_{x \in C} \sum_{i=1}^d x^{(i)} \log x^{(i)} \leq \log d\end{aligned}$$

$$\text{Hence, } f_{\text{best}}^t - f^* = O\left(\frac{L_f \sqrt{D_{\max}}}{\sqrt{\rho}} \cdot \frac{\log t}{\sqrt{t}}\right) = O\left(L_{f,\infty} \cdot \sqrt{\log d} \cdot \frac{\log t}{\sqrt{t}}\right)$$

Comparison:

	Euclidean	KL Divergence
Convergence Rate	$O(L_{f,2} \cdot \frac{\log t}{\sqrt{t}})$	$O(L_{f,\infty} \cdot \sqrt{\log d} \cdot \frac{\log t}{\sqrt{t}})$

Note that $\|\nabla f\|_\infty \leq \|\nabla f\|_2 \leq \sqrt{d} \cdot \|\nabla f\|_\infty$ (why?)

Hence,

$$\frac{1}{\sqrt{d}} \leq \frac{L_{f,\infty}}{L_{f,2}} \leq 1$$

MD with KL divergence has a better convergence rate!

Newton's Method

Recall: Non-Homogeneity and “Scaled” Gradients

$$\underset{x \in \mathbb{R}^2}{\text{minimize}} \quad f(x) := \frac{1}{2}(x - x^*)^\top Q(x - x^*)$$

Q is P.d.

- Suppose $Q = [Q_{11}, 0; 0, Q_{22}]$ is a diagonal matrix with $Q_{11} \gg Q_{22}$

- Idea: Accelerate GD by *scaling the gradient*

$$x_{t+1} = x_t - \eta_t Q^{-1} \nabla f(x_t) = x_t - \eta_t \cancel{Q} \cancel{(\nabla f(x_t))} = x_t - \eta_t (x_t - x^*)$$

Notably, in quadratic problems, we also have

$$\nabla^2 f(x) = \cancel{Q} \Leftrightarrow \cancel{Q}^\top = (\nabla^2 f(x))^{-1}$$

- Question: How about “scaled gradients” for objectives *beyond quadratic functions*?

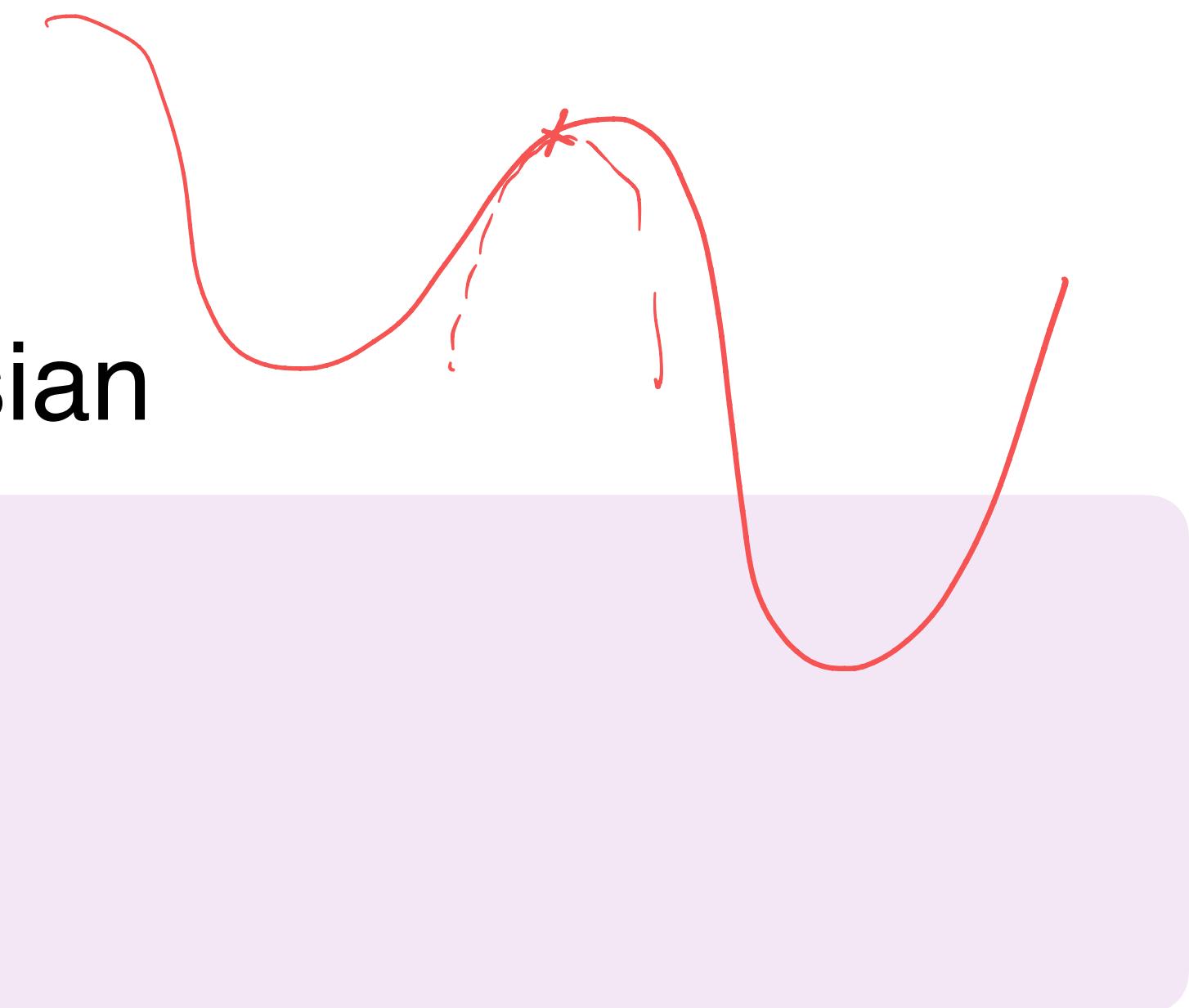
Newton's Method

For an unconstrained problem: $\underset{x \in \mathbb{R}^d}{\text{minimize}} f(x)$

Newton's method = “scaled” gradient by inverse of Hessian

$$x_{t+1} = x_t - \eta_t (\nabla^2 f(x_t))^{-1} \nabla f(x_t)$$

$=: \Delta x_{NT}$ Newton's step



- ▶ Pure Newton's step: Set $\eta_t = 1$
- ▶ Damped Newton's step: η_t is determined by “backtracking line search” (discussed later)

-
- ▶ Question: Is Newton's method a “descent” method?

$$\nabla f(x_t)^T \left(\left(\nabla^2 f(x_t) \right)^+ \nabla f(x_t) \right) \geq 0$$

$$\nabla f(x)^T z \leq 0$$

If you choose an update direction z ,

Interpretations of Newton's Step: (1) Minimizer of Second-Order Approximation

Newton's step

$$\Delta x_{NT} := (\nabla^2 f(x))^{-1} \nabla f(x)$$

FoNC:

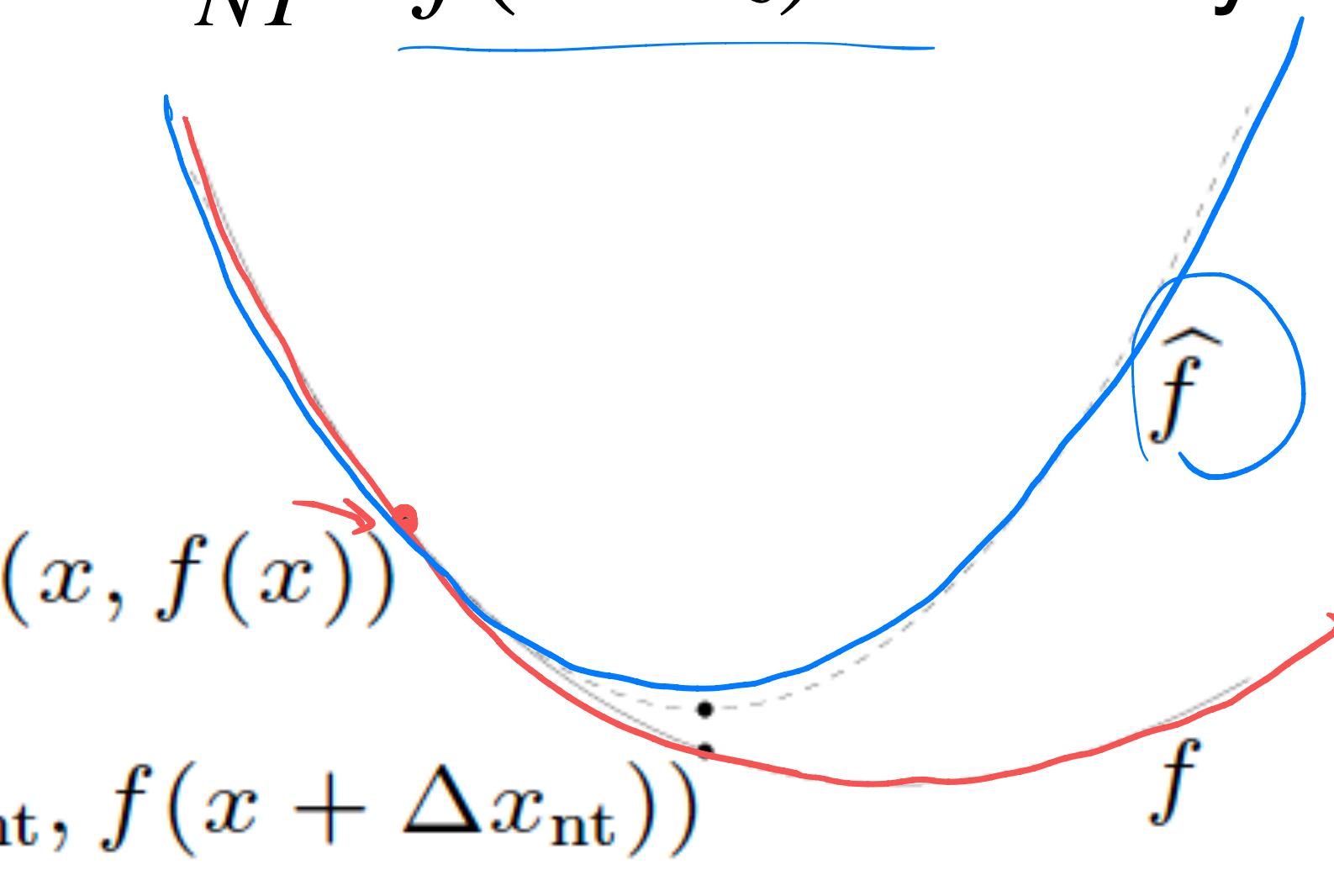
$$\begin{aligned}\nabla f(x) &= \nabla^2 f(x) \cdot z \\ \Rightarrow z &= (\nabla^2 f(x))^{-1} \nabla f(x)\end{aligned}$$

- Second-order approximation of f at some given point x :

$$\hat{f}(x + z) = f(x) + \nabla f(x)^T z + \frac{1}{2} z^T \nabla^2 f(x) z$$

- $\hat{f}(x + z)$, as a function of z , is minimized at $z = \Delta x_{NT}$ if $\hat{f}(x + z)$ is strictly convex in z
- If $\nabla f(x) = 0$, then $\Delta x_{NT} = 0$

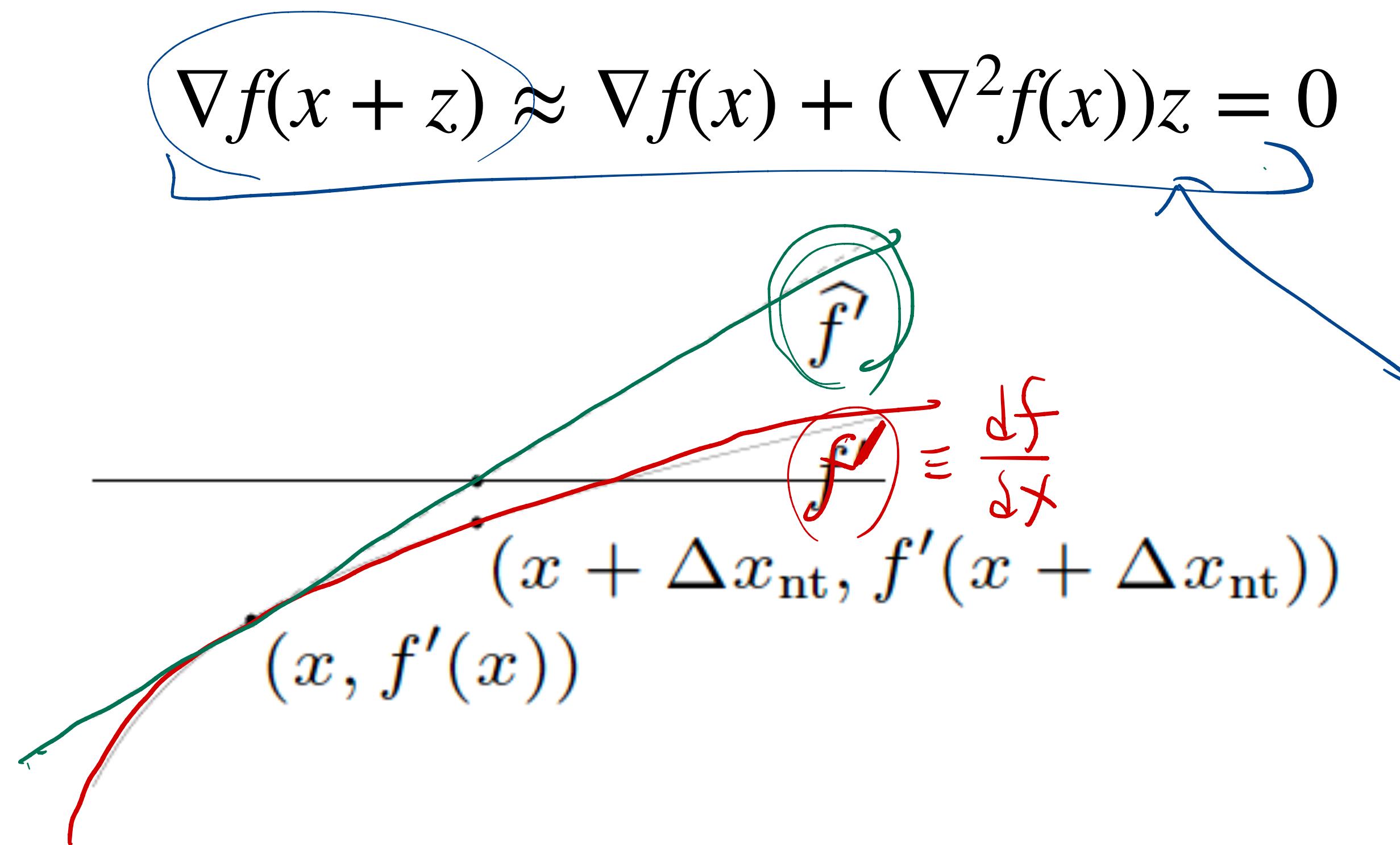
Takeaway: Newton's Step finds the minimizer of second-order approximation $(x, f(x))$ (under f strictly convex) $(x + \Delta x_{nt}, f(x + \Delta x_{nt}))$



Interpretations of Newton's Step: (2) Solution to Linearized Optimality Condition

Newton's step $\Delta x_{NT} := (\nabla^2 f(x))^{-1} \nabla f(x)$

- First-order necessary condition near x under linear approximation



$$\hat{f}(x+z) = f(x) + \nabla f(x)^T z + \frac{1}{2} z^T \nabla^2 f(x) z$$
$$\nabla \hat{f}(x+z) = \nabla f(x) + (\nabla^2 f(x)) \cdot z$$

Features of Newton's Method

1. Invariance under (invertible) linear transformation of coordinates

- Let $P : \mathbb{R}^d \rightarrow \mathbb{R}^d$ be an invertible linear transformation and $x = \underline{P}y$
- Then, we have $\underline{x} + \Delta x_{NT} = P(y + \Delta y_{NT})$

$$\xrightarrow{\quad\quad\quad} \quad\quad\quad \xleftarrow{\quad\quad\quad}$$

$$\circlearrowleft \quad \circlearrowright$$

$$\|x_{t+1} - x^*\| \leq \gamma \cdot \|x_t - x^*\|^2$$

2. "Superlinear" "local" convergence

- Superlinear: $\|x_{t+1} - x^*\| \leq \gamma \|x_t - x^*\|^\alpha$ with $\alpha > 1$ (x^* is a stationary point)
- Local: The above holds only if one starts within a small neighborhood of x^*
- Note: The above does NOT imply convergence to local minimum!

$$\swarrow \quad \nwarrow$$

$$\circlearrowleft \quad \circlearrowright$$

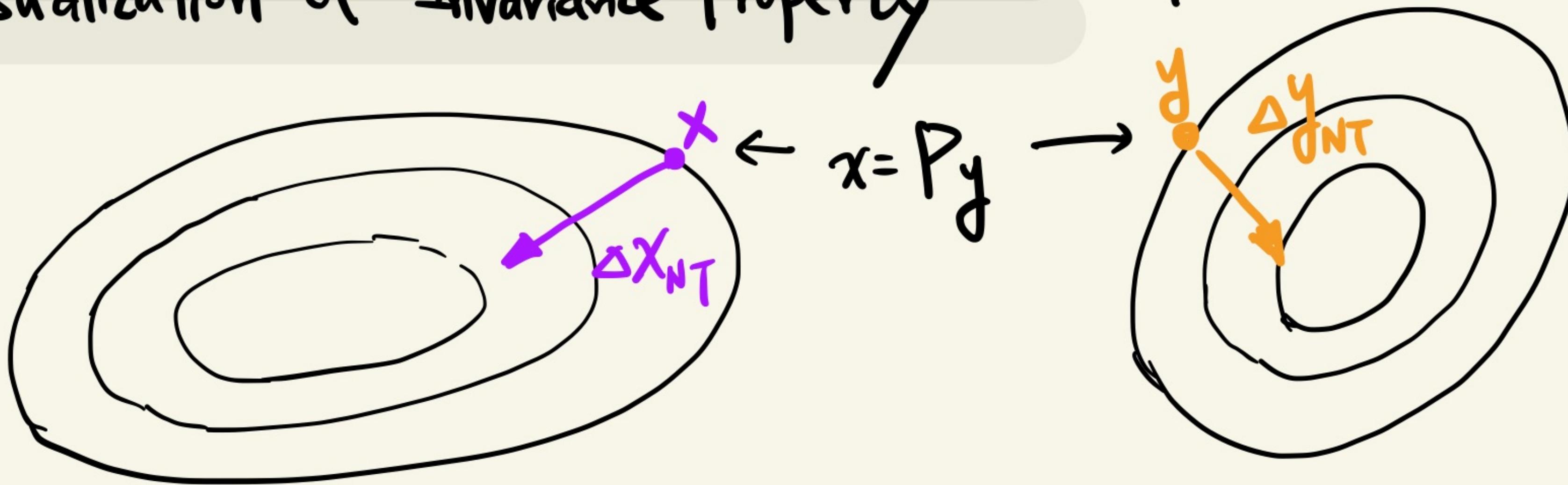
Recall: Terminology of Convergence Rates

- Let $e(x)$ denote the distance from optimality
 - Example: $e(x) = \|x - x^*\|$
 - Example: $e(x) = |f(x) - f(x^*)|$
- **Rate of convergence:** The limit of the ratio of successive errors

$$\lim_{k \rightarrow \infty} \frac{e(x_{k+1})}{e(x_k)} = \beta$$

- If $\beta = 1$: We call it a **sub-linear rate** of convergence
- If $\beta \in (0,1)$: We call it a **linear rate** of convergence
- If $\beta = 0$: We call it a **super-linear rate** of convergence

Visualization of Invariance Property

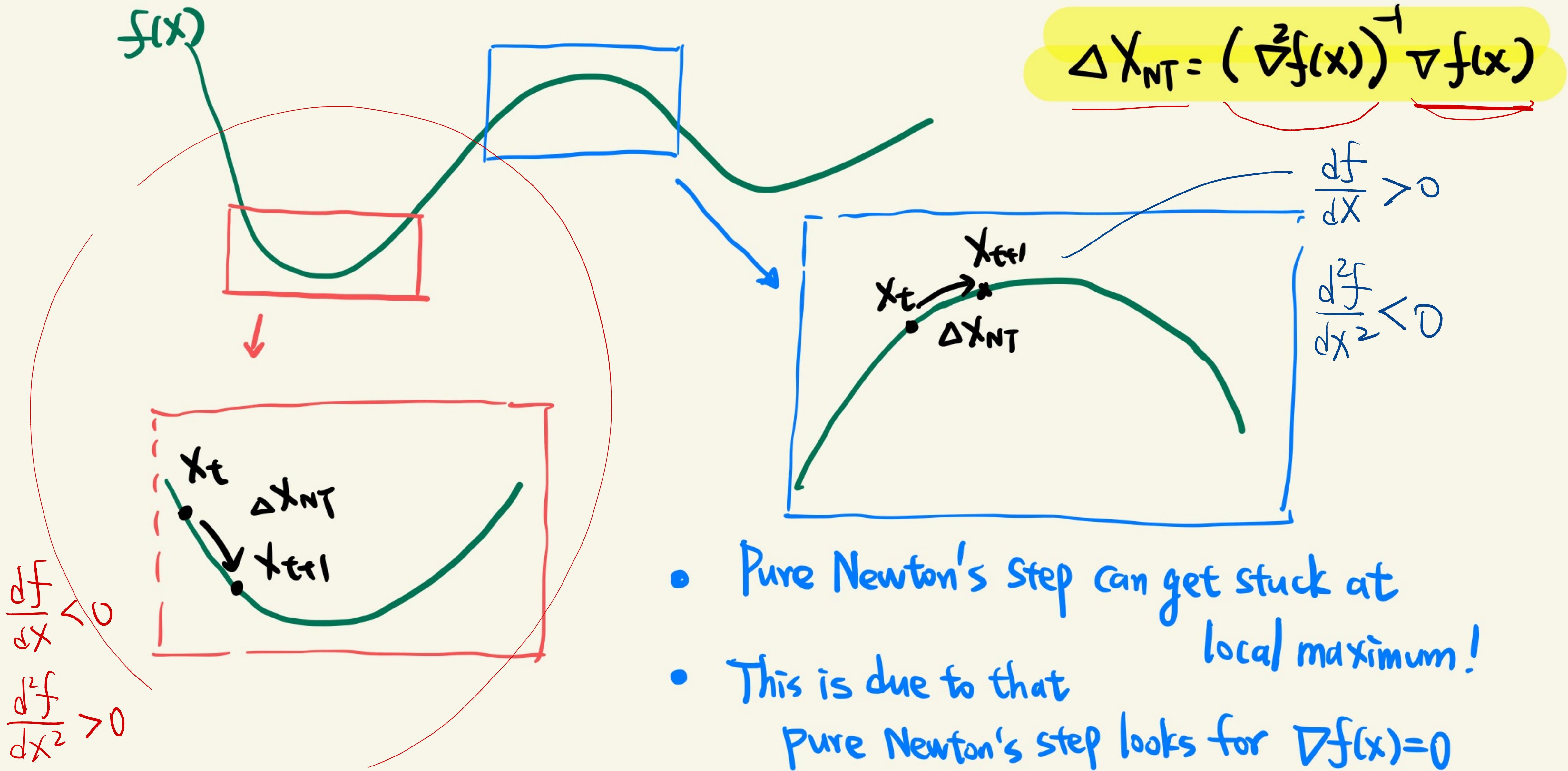


- Define $h(y) = f(Py)$

$$\Rightarrow \begin{cases} \nabla h(y) = P^T \nabla f(Py) = P^T \nabla f(x) \\ \nabla^2 h(y) = P^T (\nabla^2 f(x)) P \end{cases}$$

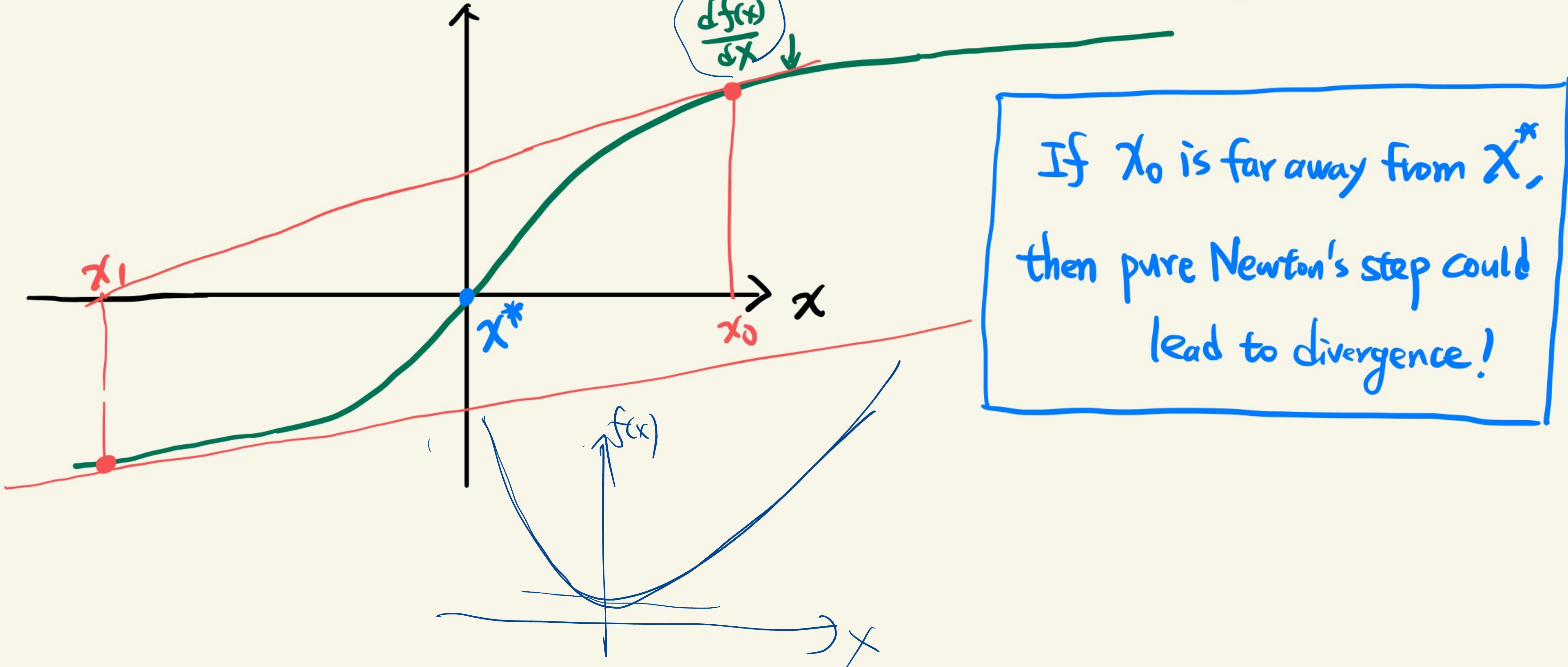
$$\begin{aligned} \Delta y_{NT} &= -\left(P^T \nabla^2 f(x) P \right)^{-1} (P^T \nabla f(x)) \\ &= P^{-1} \left(-(\nabla^2 f(x))^{-1} \nabla f(x) \right) \end{aligned}$$

Convergence Issue of Pure Newton's Step



Convergence Issue of Pure Newton's Step

Consider a "strongly-convex" function $f(x) : \mathbb{R}^l \rightarrow \mathbb{R}^l$ with $\frac{\partial f(x)}{\partial x}$ as :

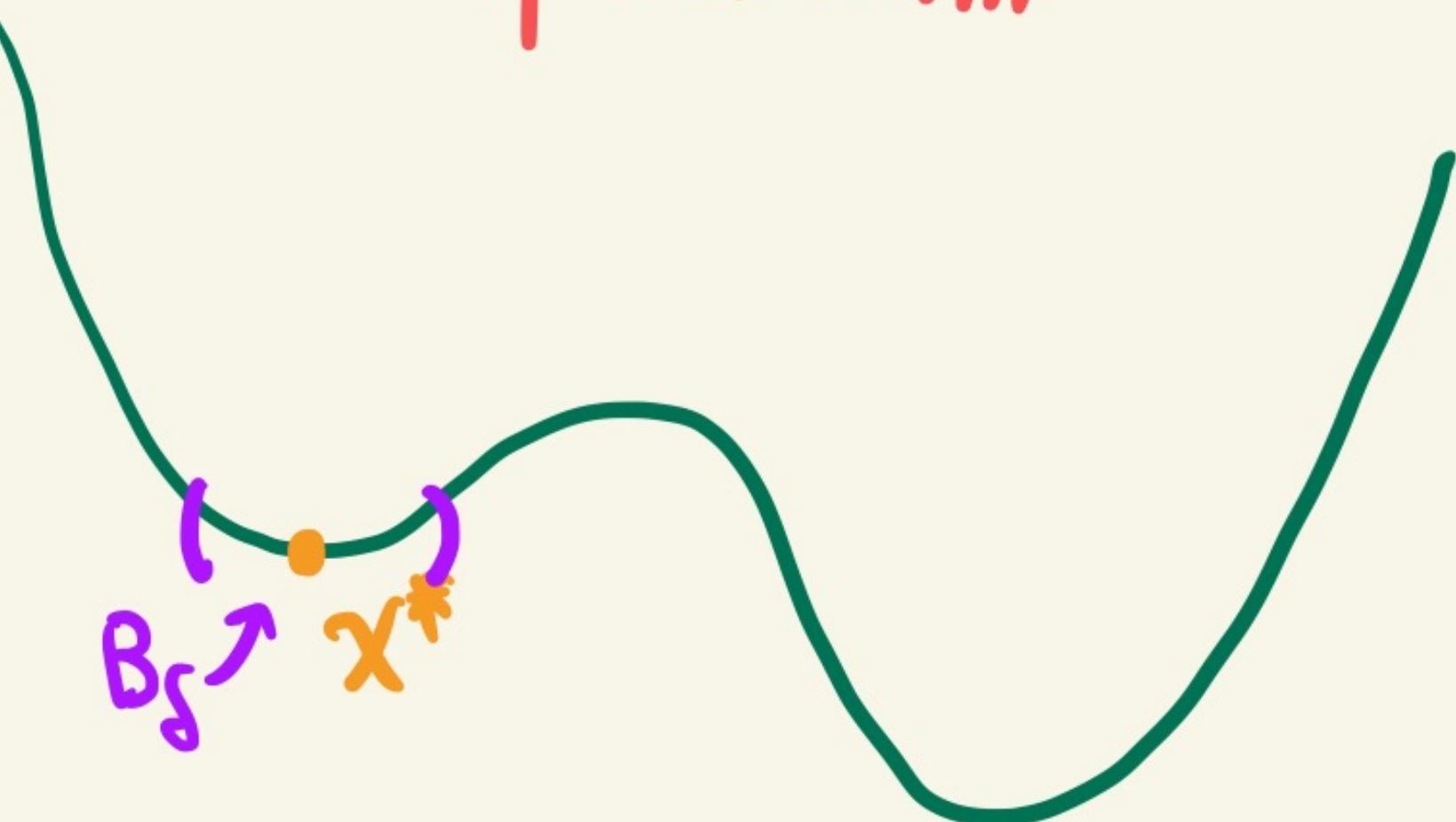


Convergence of Newton's Method = Local Convergence

Theorem Let $f(x)$ be twice continuously differentiable (could be non-convex)

- Let x^* be a stationary point (i.e., $\nabla f(x^*) = 0$) and define $B_\delta := \{x' : \|x' - x^*\| \leq \delta\}$
- Hessian of $f(x)$ is L -Lipschitz continuous, i.e., $\|\nabla^2 f(x) - \nabla^2 f(y)\| \leq L \cdot \|x - y\|$
- Suppose there exists $\delta > 0$ s.t. $(\nabla^2 f(x))^{-1}$ exists and $\|\nabla^2 f(x)^{-1}\| \leq M$, for all $x \in B_\delta$.
↑ spectral norm
- The initial point $x_0 \in B_\delta$.

Then, $\|x_{t+1} - x^*\| \leq \frac{LM}{2} \cdot \|x_t - x^*\|^2$



A Useful Lemma: Gradient Theorem for Line Integrals

Lemma: Let $g(x): \mathbb{R}^d \rightarrow \mathbb{R}$ be a scalar-valued differentiable function.

Then, for any "continuous curve" ϕ which starts at x_1 and ends at x_2

$$\int_{\phi} \nabla g(z)^T dz = g(x_2) - g(x_1)$$

- If ϕ is a line segment from x_1 to x_2 :
- $$g(x_2) - g(x_1) = \int_{\phi} \nabla g(z)^T dz = \left(\int_0^1 \nabla g(x_1 + t \cdot (x_2 - x_1))^T dt \right) (x_2 - x_1)$$

Proof of Convergence

$$\underline{\text{Step 1:}} \quad \|x_{t+1} - x^*\| = \|x_t - x^* - (\nabla^2 f(x_t))^{-1} \cdot \nabla f(x_t)\|$$

$$= \|(\nabla^2 f(x_t))^{-1} \left((\nabla^2 f(x_t)) \cdot (x_t - x^*) - \underbrace{\nabla f(x_t)}_{\text{"}} \right)\|$$

$$= \|(\nabla^2 f(x_t))^{-1} \left(\int_0^1 (\nabla^2 f(x_t)) - \nabla^2 f(x^* + s \cdot (x_t - x^*))^T ds \right) \cdot (x_t - x^*)\|$$

$$\leq \|(\nabla^2 f(x_t))^{-1}\| \cdot \left\| \int_0^1 (\nabla^2 f(x_t)) - \nabla^2 f(x^* + s \cdot (x_t - x^*))^T ds \right\| \cdot \|x_t - x^*\|$$

(Cont.).

Step 2: $\left\| \int_0^1 (\nabla^2 f(x_t)) - \nabla^2 f(x^* - s(x_t - x^*)) ds \right\|$

Triangle inequality \downarrow

$$\leq \int_0^1 \left\| \nabla^2 f(x_t) - \nabla^2 f(x^* - s(x_t - x^*)) \right\| ds$$
$$\leq L \cdot s \cdot \|x_t - x^*\| \quad \text{by smoothness}$$
$$\leq \frac{L}{2} \|x_t - x^*\|$$

Step 3: By Combining Step 1 and Step 2,

$$\|x_{t+1} - x^*\| \leq \frac{LM}{2} \cdot \|x_t - x^*\|^2.$$

□

Remark on Newton's Decrement

$$\lambda(x)^2 := \nabla f(x)^\top (\nabla^2 f(x))^{-1} \nabla f(x) \equiv \nabla f(x)^\top \Delta x_{NT}$$

- ▶ Newton's decrement = $f(x) - \inf_z \hat{f}(z)$ (where $\hat{f}(z)$ is the second-order approximation)

$$f(x) - \inf_z \hat{f}(z) = f(x) - \hat{f}(x + \Delta x_{NT}) = \frac{1}{2} \lambda(x)^2$$

- ▶ Hence, Newton's decrement is a good estimate of sub-optimality gap $f(x) - f(x^*)$

Summary: Pros and Cons of Newton's Method

$$\text{minimize}_{x \in \mathbb{R}^d} \quad f(x)$$

$$x_{t+1} = x_t - \eta_t \underbrace{(\nabla^2 f(x_t))^{-1}}_{=: \Delta x_{NT}} \nabla f(x_t)$$

-
- ▶ Pros
 - ▶ Fast convergence (cf. iteration complexity = $O(\log \log \frac{1}{\epsilon})$)
 - ▶ Cons
 - ▶ Requires storing and inverting Hessian $\nabla^2 f(x) \in R^{d \times d}$
 - ▶ One single iteration could take forever
 - ▶ Prohibitively large storage overhead

Next Topic: Can we do Newton's step without Hessian?

Quasi-Newton Methods!

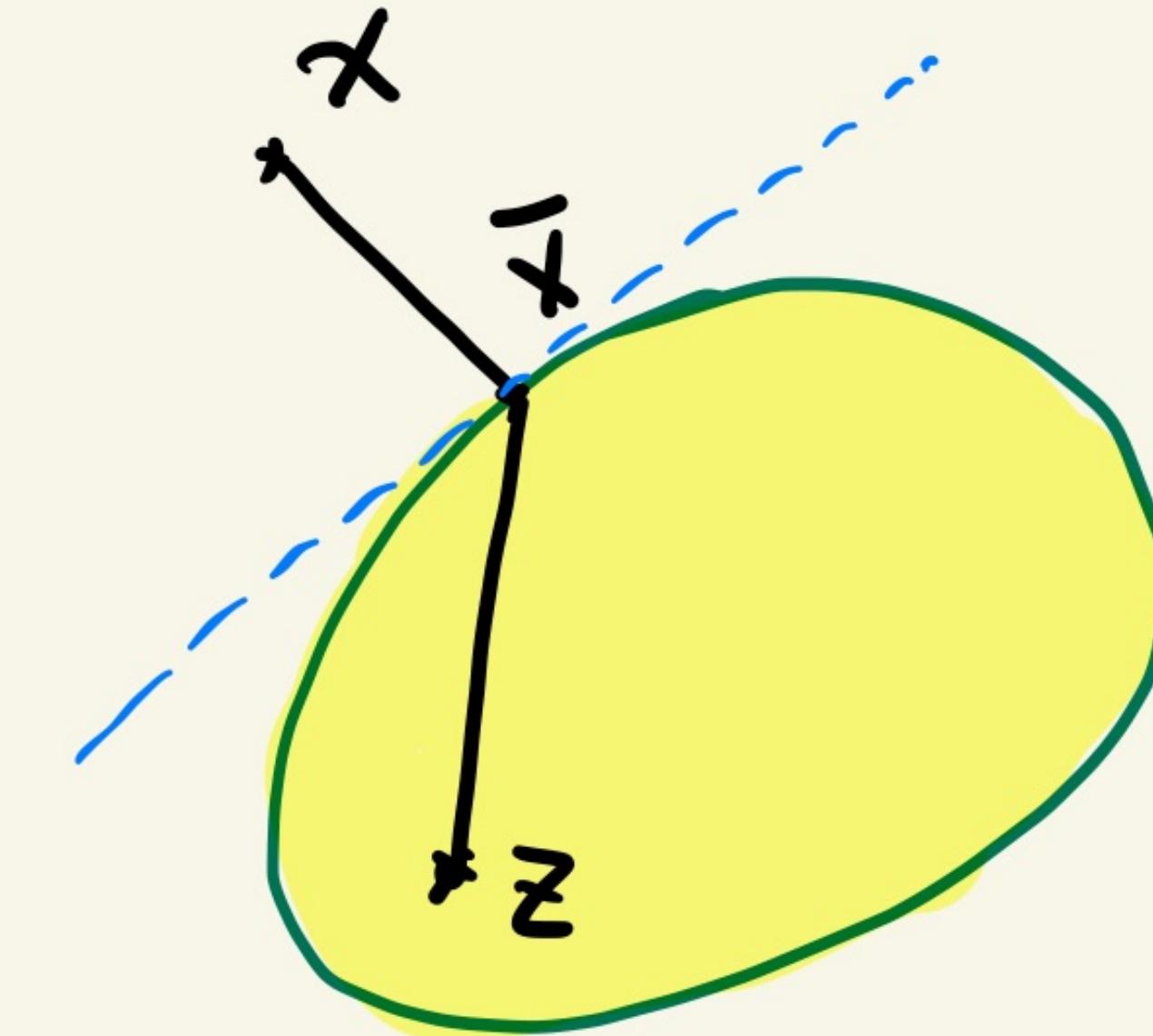
Appendix: Proof of Convergence of Mirror Descent

Amir Beck and Marc Teboulle, “Mirror descent and nonlinear projected subgradient methods for convex optimization,”
Operations Research Letters, 2003.

Generalized Pythagorean Theorem

Euclidean Case

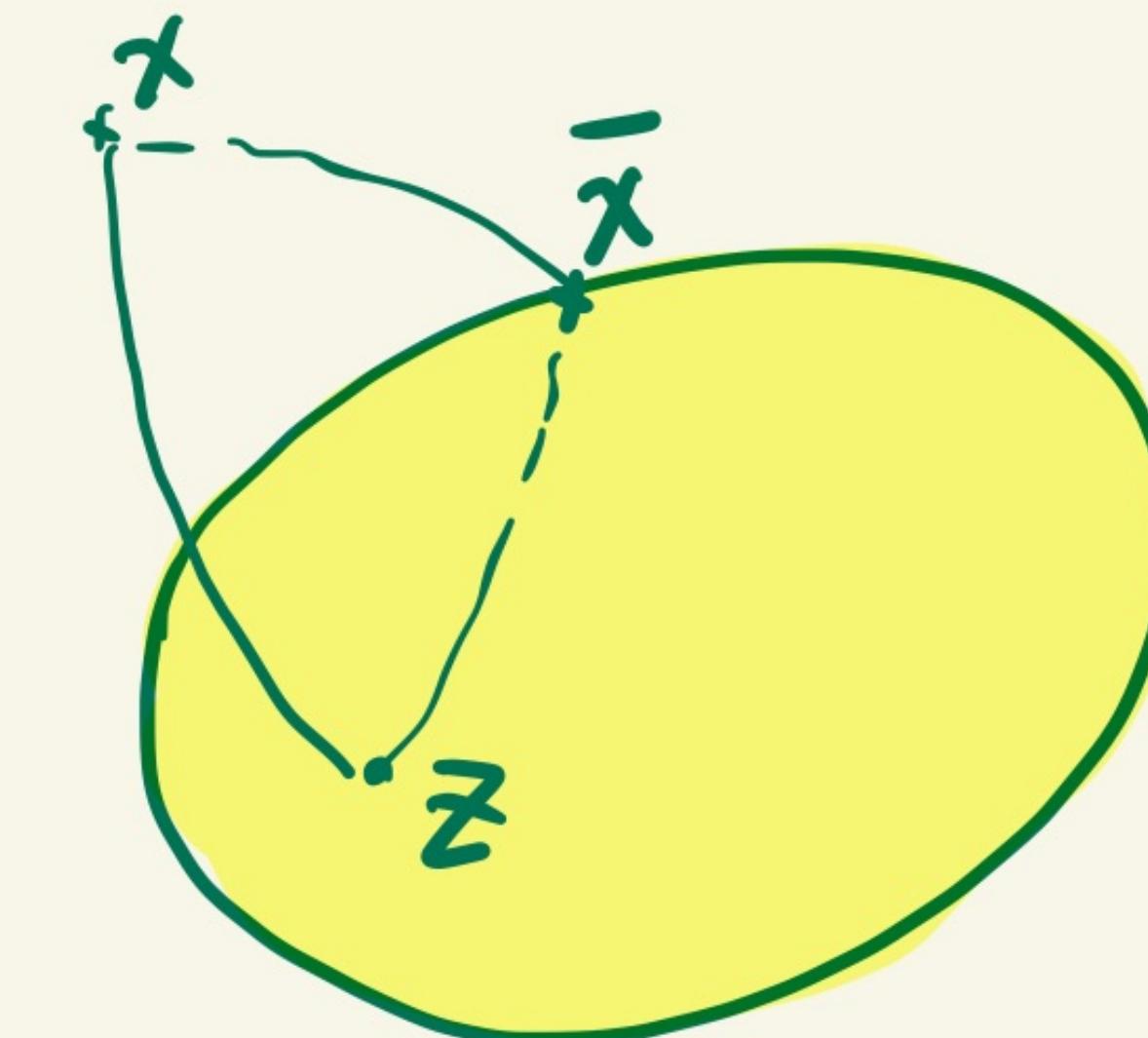
$$\|z-x\|^2 \geq \|\bar{x}-z\|^2 + \|x-\bar{x}\|^2$$



Bregman Divergence

$$D_\phi(z||x) \geq D_\phi(z||\bar{x}) + D_\phi(\bar{x}||x).$$

(Proof: HW2)



Hölder's Inequality

Lemma: Let $p, q \in [1, \infty]$ so that $\frac{1}{p} + \frac{1}{q} = 1$. Then, for any two vectors

$x = (x_1, \dots, x_n)$, $y = (y_1, \dots, y_n)$, we have

$$\left| \sum_{i=1}^n x_i y_i \right| \leq \|x\|_p \cdot \|y\|_q.$$

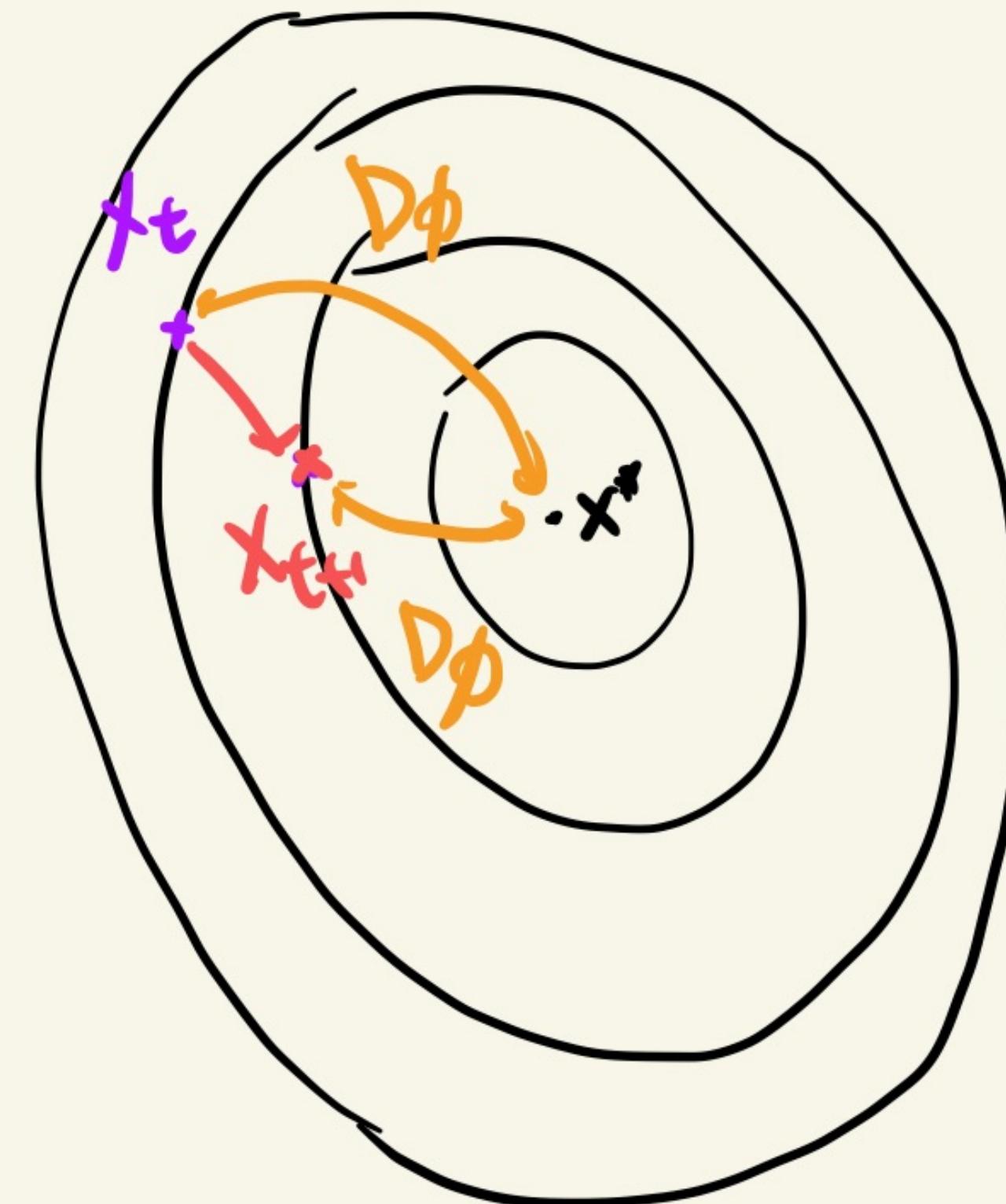
Remark: If $p=2, q=2$, Hölder's inequality is essentially Cauchy-Schwarz

A Fundamental Inequality of MD

Lemma

$$\gamma_t \cdot (f(x_t) - f^*) \leq D_\phi(x^* \| x_t) - D_\phi(x^* \| x_{t+1}) + \frac{\gamma_t^2 L_f^2}{2 \cdot \rho}$$

Question: What's the intuition?



Proof of Lemma :

Step 1: $f(x_t) - f(x^*) \leq \nabla f(x_t)^T (x_t - x^*) \dots \text{ (By convexity)}$

$$= \frac{1}{\gamma_t} \cdot (\nabla \phi(x_t) - \nabla \phi(\tilde{x}_{t+1}))^T (x_t - x^*) \dots \text{ (By the "mirror map")}$$

$$= \frac{1}{\gamma_t} \cdot (D_\phi(x^* \| x_t) + D_\phi(x_t \| \tilde{x}_{t+1}) - D_\phi(x^* \| \tilde{x}_{t+1})) \dots \text{ (By Three-Point Lemma)}$$

$$\leq \frac{1}{\gamma_t} \cdot \underbrace{(D_\phi(x^* \| x_t) + D_\phi(x_t \| \tilde{x}_{t+1}))}_{\text{green bracket}} - \underbrace{(D_\phi(x^* \| x_{t+1}) + D_\phi(x_{t+1} \| \tilde{x}_{t+1}))}_{\text{purple bracket}} \dots \text{ (By Generalized Pythagorean Theorem)}$$

(Cont.) .

Step 2: It is sufficient to show the following claim:

$$D_\phi(x_t, \tilde{x}_{t+1}) - D_\phi(x_{t+1}, \tilde{x}_{t+1}) \leq \frac{(l_t \cdot L_f)^2}{2 \cdot \rho} \quad (*)$$

(which would naturally lead to our required lemma)

$$\begin{aligned} & D_\phi(x_t, \tilde{x}_{t+1}) - D_\phi(x_{t+1}, \tilde{x}_{t+1}) \\ &= \phi(x_t) - \phi(x_{t+1}) - \nabla \phi(\tilde{x}_{t+1})^T (x_t - x_{t+1}) \dots \text{ (By Bregman divergence)} \\ &\leq \nabla \phi(x_t)^T (x_t - x_{t+1}) - \frac{\rho}{2} \|x_t - x_{t+1}\|^2 - \nabla \phi(\tilde{x}_{t+1})^T (x_t - x_{t+1}) \\ &\quad \dots (\phi \text{ is } \rho\text{-strongly convex}) \end{aligned}$$

(cont.).

$$= (\nabla \phi(x_t) - \nabla \phi(\tilde{x}_{t+1}))^\top (x_t - x_{t+1}) - \frac{\rho}{2} \|x_t - x_{t+1}\|^2$$

$$= \eta_t \cdot \nabla f(x_t)^\top (x_t - x_{t+1}) - \frac{\rho}{2} \|x_t - x_{t+1}\|^2 \quad \dots \text{ (By mirror descent update)}$$

$$\leq \eta_t \cdot L_f \cdot \|x_t - x_{t+1}\| - \frac{\rho}{2} \|x_t - x_{t+1}\|^2 \quad \dots \text{ (By "Holder's inequality")}$$

$$\leq \frac{(\eta_t \cdot L_f)^2}{2\rho} \quad \dots \text{ (By completing the square)}$$

Proof of Main Theorem

Step 1: By taking the telescoping sum of $\sum_{\tau} (f(x_\tau) - f^*)$, we have

$$\sum_{\tau=0}^t \gamma_\tau (f(x_\tau) - f^*) \leq D_\phi(x^* \| x_0) - D_\phi(x^* \| x_{t+1}) + \frac{L_f^2 \cdot \sum_{\tau=0}^t \gamma_\tau^2}{2 \cdot \rho}$$

$$\leq \sup_{x \in C} D_\phi(x \| x_0) + \frac{L_f^2 \cdot \sum_{\tau=0}^t \gamma_\tau^2}{2 \cdot \rho}$$

Step 2: Moreover, we have

$$f_{\text{best}}^t - f^* \leq \frac{\sum_{\tau=0}^t \gamma_\tau (f(x_\tau) - f^*)}{\sum_{\tau=0}^t \gamma_\tau}$$