

Definitions	
<ul style="list-style-type: none"> <li><b>Differentiable:</b> <math>f : \mathbb{R}^d \rightarrow \mathbb{R}</math> is differentiable if <math display="block">f(\mathbf{y}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{r(\mathbf{x} - \mathbf{y})}{\ \mathbf{x} - \mathbf{y}\ },</math> where <math>\lim_{\mathbf{v} \rightarrow \mathbf{0}} \frac{\ r(\mathbf{v})\ }{\ \mathbf{v}\ } = 0</math>. </li> <li><b>Spectral norm:</b> <math>\ A\ _2 = \sup_{\ \mathbf{x}\ =1} \ A\mathbf{x}\ </math> (largest eigenvalue).</li> <li><b>Positive semi-definite:</b> <math>\forall \mathbf{x} \in \mathbb{R}^d: \mathbf{x}^\top A \mathbf{x} \geq 0</math>.</li> <li><b>B-Lipschitz:</b> <math>\ f(\mathbf{x}) - f(\mathbf{y})\  \leq B\ \mathbf{x} - \mathbf{y}\  \Leftrightarrow \ \nabla f(\mathbf{x})\  \leq B</math>.</li> <li><b>Convex set:</b> <math>\forall \mathbf{x}, \mathbf{y} \in X, \lambda \in [0, 1]: \lambda \mathbf{x} + (1 - \lambda)\mathbf{y} \in X</math>.</li> <li><b>Convexity:</b> <math>\forall \mathbf{x}, \mathbf{y} \in \text{dom}(f)</math> and <math>\forall \lambda \in [0, 1]</math>, <ol style="list-style-type: none"> <li><math>f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y})</math>.</li> <li><math>f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle</math>.</li> <li><math>\langle \nabla f(\mathbf{x}) + \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq 0</math>.</li> <li><math>\nabla^2 f(\mathbf{x})</math> is positive semi-definite.</li> </ol> </li> <li><b>L-smoothness:</b> <math>\forall \mathbf{x}, \mathbf{y} \in \text{dom}(f)</math>, <ol style="list-style-type: none"> <li><math>\ \nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\  \leq L\ \mathbf{x} - \mathbf{y}\ </math>.</li> <li><math>g(\mathbf{x}) := \frac{L}{2}\ \mathbf{x}\ ^2 - f(\mathbf{x})</math> is convex.</li> <li><math>f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2}\ \mathbf{x} - \mathbf{y}\ ^2</math> (canonical).</li> <li><math>\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \leq L\ \mathbf{x} - \mathbf{y}\ ^2</math>.</li> <li><math>\ \nabla^2 f(\mathbf{x})\ _2 \leq L</math>.</li> <li>Coordinate-wise: <math>f(\mathbf{x} + \lambda \mathbf{e}_i) \leq f(\mathbf{x}) + \lambda \nabla_i f(\mathbf{x}) + \frac{L_i}{2} \lambda^2, \forall \lambda \in \mathbb{R}</math>. Relations: [5] <math>\Leftrightarrow</math> [1] <math>\Rightarrow</math> [2] <math>\Leftrightarrow</math> [3] <math>\Leftrightarrow</math> [4] (If convex, all <math>\Leftrightarrow</math>).</li> </ol> </li> <li><b><math>\mu</math>-strong convexity:</b> <math>\forall \mathbf{x}, \mathbf{y} \in \text{dom}(f)</math>, <ol style="list-style-type: none"> <li><math>f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\mu}{2}\ \mathbf{x} - \mathbf{y}\ ^2</math>.</li> <li><math>g(\mathbf{x}) := f(\mathbf{x}) - \frac{\mu}{2}\ \mathbf{x}\ ^2</math> is convex.</li> <li><math>\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq \mu\ \mathbf{x} - \mathbf{y}\ ^2</math> (needs proof).</li> <li><math>\mu</math>-SC <math>\Rightarrow</math> PL inequality: <math>\frac{1}{2}\ \nabla f(\mathbf{x})\ ^2 \geq \mu(f(\mathbf{x}) - f^*)</math>.</li> </ol> </li> <li><b>Subgradient:</b> <math>\mathbf{g} \in \partial f(\mathbf{x}) \Leftrightarrow f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{g}, \mathbf{y} - \mathbf{x} \rangle, \forall \mathbf{y} \in \text{dom}(f)</math>.</li> <li><b>Conjugate function:</b> <math>f^*(\mathbf{y}) := \sup_{\mathbf{x} \in \text{dom}(f)} \langle \mathbf{x}, \mathbf{y} \rangle - f(\mathbf{x})</math>.</li> </ul>	<ul style="list-style-type: none"> <li><math>\left\  \int_0^1 \nabla h(t) dt \right\  \leq \int_0^1 \ \nabla h(t)\  dt</math>.</li> <li><math>\int_0^1 c dt = c, \quad \int_0^1 t dt = \frac{1}{2}</math>.</li> <li><b>Subgradient calculus:</b> <ol style="list-style-type: none"> <li><math>h(\mathbf{x}) = \alpha f(\mathbf{x}) + \beta g(\mathbf{x}) \Rightarrow \partial h(\mathbf{x}) = \alpha \cdot \partial f(\mathbf{x}) + \beta \cdot g(\mathbf{x})</math>.</li> <li><math>h(\mathbf{x}) = f(A\mathbf{x} + \mathbf{b}) \Rightarrow \partial h(\mathbf{x}) = A^\top \partial f(A\mathbf{x} + \mathbf{b})</math>.</li> <li><math>h(\mathbf{x}) = \max f_i(\mathbf{x}) \Rightarrow \partial h(\mathbf{x}) = \text{conv}(\{\partial f_i(\mathbf{x}) \mid f_i(\mathbf{x}) = h(\mathbf{x})\})</math>.</li> </ol> </li> </ul>
Optimality lemmas (assume convexity)	
<p>The constrained and non-differentiable cases are useful when the update rule contains an argmin.</p> <ul style="list-style-type: none"> <li><math>\mathbf{x}^*</math> is a local minimum.</li> <li><math>\nabla f(\mathbf{x}^*) = \mathbf{0}</math>.</li> <li>Constrained: <math>\nabla f(\mathbf{x}^*)^\top (\mathbf{x} - \mathbf{x}^*) \geq 0, \forall \mathbf{x} \in X</math>.</li> <li>Non-differentiable: <math>\mathbf{0} \in \partial f(\mathbf{x}^*)</math>.</li> </ul>	
Common tricks	
<ul style="list-style-type: none"> <li><b>Rearrange the update rule</b> for an equality—e.g., <math>\nabla f(\mathbf{x}_t) = \frac{\mathbf{x}_t - \mathbf{x}_{t+1}}{\gamma_t}</math>.</li> <li>Define <math>h(t) := f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))</math>, where <math>h'(t) = \nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))^\top (\mathbf{y} - \mathbf{x})</math> and use with fundamental theorem of calculus, <math display="block">f(\mathbf{y}) - f(\mathbf{x}) = \int_0^1 \nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))^\top (\mathbf{y} - \mathbf{x}) dt.</math> Or, mean-value theorem, <math display="block">\nabla f(\mathbf{x} + c(\mathbf{y} - \mathbf{x}))^\top (\mathbf{y} - \mathbf{x}) = f(\mathbf{y}) - f(\mathbf{x}), \quad \exists c \in (0, 1).</math> </li> <li>Projection is <b>non-expansive</b>: <math>\ \Pi_X(\mathbf{x}) - \Pi_X(\mathbf{y})\  \leq \ \mathbf{x} - \mathbf{y}\ </math>.</li> <li><math>\min_{1 \leq t \leq T} f(\mathbf{x}_t) - f^* \leq \frac{\sum_{t=1}^T \gamma_t (f(\mathbf{x}_t) - f^*)}{\sum_{t=1}^T \gamma_t}</math>.</li> <li><b>Telescoping sum</b> inequality: <math>\sum_{t=1}^T \ \mathbf{x}_t - \mathbf{x}^*\ ^2 - \ \mathbf{x}_{t+1} - \mathbf{x}^*\ ^2 \leq \ \mathbf{x}_1 - \mathbf{x}^*\ ^2</math>.</li> <li><math>f^* \leq f(\mathbf{x}), \forall \mathbf{x} \in X</math> can sometimes be useful to bound <math>f(\mathbf{x}_t) - f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - f^*</math>.</li> <li><math>\max\{a, b\} \leq a + b</math> if <math>a, b \geq 0</math>.</li> </ul>	
Expectation and variance for SGD	
<ul style="list-style-type: none"> <li><math>\text{Var}[X] := \mathbb{E}[(X - \mathbb{E}[X])^2]</math></li> <li><math>\text{Var}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2</math>  <math>\Rightarrow \mathbb{E}\ \nabla f(\mathbf{x}_t, \xi_t)\ ^2 = \ \nabla F(\mathbf{x}_t)\ ^2 + \mathbb{E}\ \nabla f(\mathbf{x}_t, \xi_t) - \nabla F(\mathbf{x}_t)\ ^2 \leq \ \nabla F(\mathbf{x}_t)\ ^2 + \sigma^2</math>.</li> <li><b>Law of total expectation:</b> <math>\mathbb{E}[X] = \mathbb{E}_Y[\mathbb{E}_X[X \mid Y]]</math>.</li> <li><b>Law of total var.:</b> <math>\text{Var}[Y] = \mathbb{E}_X[\text{Var}_Y[Y \mid X]] + \text{Var}_Y[\mathbb{E}_X[Y \mid X]]</math>.</li> <li><math>\text{Var}[X - Y] = \text{Var}[X] + \text{Var}[Y] - 2 \cdot \text{Cov}(X, Y)</math>.</li> <li><math>\text{Var}[\alpha X] = \alpha^2 \text{Var}[X], \text{Var}[X + \beta] = \text{Var}[X]</math>.</li> </ul>	
Risk minimization	
Non-linear programming	
Gradient descent	
<ul style="list-style-type: none"> <li><b>Update rule:</b> <math>\mathbf{x}_{t+1} = \mathbf{x}_t - \gamma \nabla f(\mathbf{x})</math>.</li> <li><b>VA:</b> <math>\sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f^*) \leq \frac{\gamma}{2} \sum_{t=0}^{T-1} \ \nabla f(\mathbf{x}_t)\ ^2 + \frac{1}{2\gamma} \ \mathbf{x}_0 - \mathbf{x}^*\ ^2</math>.</li> <li><b>1st-order convexity</b> on <math>(\mathbf{x}^*, \mathbf{x}_t) \Rightarrow \nabla f(\mathbf{x}_t) = \frac{\mathbf{x}_t - \mathbf{x}_{t+1}}{\gamma} \Rightarrow</math> Cosine theorem <math>\Rightarrow \mathbf{x}_t - \mathbf{x}_{t+1} = \gamma \nabla f(\mathbf{x}_t) \Rightarrow</math> Telescoping sum.</li> <li><b>Sufficient decrease:</b> <math>f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L} \ \nabla f(\mathbf{x}_t)\ ^2</math>.</li> <li><b>Smoothness</b> on <math>(\mathbf{x}_{t+1}, \mathbf{x}_t) \Rightarrow \mathbf{x}_{t+1} - \mathbf{x}_t = -\frac{1}{L} \nabla f(\mathbf{x}_t)</math>.</li> <li><b>Convergence results:</b> (<math>\ \mathbf{x}_0 - \mathbf{x}^*\  \leq R</math>) <ul style="list-style-type: none"> <li>(B-Lipschitz, convex, <math>\gamma := \frac{R}{B\sqrt{T}}</math>) <math>\frac{1}{T} \sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f^*) \leq \frac{RB}{\sqrt{T}}</math>.</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>■ Apply bounds to VA and find <math>\gamma</math> by 1st-order optimality.</li> </ul>
Lemmas	
<ul style="list-style-type: none"> <li><b>Cosine theorem:</b> All equivalent formulations, <ol style="list-style-type: none"> <li><math>\ \mathbf{x} - \mathbf{y}\ ^2 = \ \mathbf{x}\ ^2 + \ \mathbf{y}\ ^2 - 2\langle \mathbf{x}, \mathbf{y} \rangle</math>.</li> <li><math>\langle \mathbf{x}, \mathbf{y} \rangle = \frac{1}{2}(\ \mathbf{x}\ ^2 + \ \mathbf{y}\ ^2 - \ \mathbf{x} - \mathbf{y}\ ^2)</math>.</li> <li><math>\langle \mathbf{x} - \mathbf{y}, \mathbf{x} - \mathbf{z} \rangle = \frac{1}{2}(\ \mathbf{x} - \mathbf{y}\ ^2 + \ \mathbf{x} - \mathbf{z}\ ^2 - \ \mathbf{y} - \mathbf{z}\ ^2)</math>.</li> </ol> </li> <li><b>Cauchy-Schwarz:</b> <math> \langle \mathbf{x}, \mathbf{y} \rangle  \leq \ \mathbf{x}\  \ \mathbf{y}\ </math>.</li> <li><b>Hölder's inequality</b> (special case): <math> \langle \mathbf{x}, \mathbf{y} \rangle  \leq \ \mathbf{x}\ _1 \ \mathbf{y}\ _\infty</math>.</li> <li><b>Jensen's inequality</b> (<math>\varphi</math> convex, <math>a_i \geq 0</math>):  <math display="block">\varphi\left(\frac{\sum_{i=1}^m a_i \mathbf{x}_i}{\sum_{i=1}^m a_i}\right) \leq \frac{\sum_{i=1}^m a_i \varphi(\mathbf{x}_i)}{\sum_{i=1}^m a_i}.</math> </li> <li><b>Fenchel's inequality:</b> <math>\langle \mathbf{x}, \mathbf{y} \rangle \leq f(\mathbf{x}) + f^*(\mathbf{x})</math>  <math>\Rightarrow \langle \mathbf{x}, \mathbf{y} \rangle \leq \frac{1}{2}(\ \mathbf{x}\ ^2 + \ \mathbf{y}\ _*^2)</math>.</li> <li><b>Young's inequality</b> (<math>a, b \geq 0, \frac{1}{p} + \frac{1}{q} = 1</math>): <math>ab \leq \frac{a^p}{p} + \frac{b^q}{q}</math>  <math>\Rightarrow \ \mathbf{x}\  \ \mathbf{y}\  \leq \frac{1}{2}(\ \mathbf{x}\ ^2 + \ \mathbf{y}\ ^2)</math>.</li> <li><math>\frac{1}{\sqrt{d}} \ \mathbf{x}\ _2 \leq \ \mathbf{x}\ _\infty \leq \ \mathbf{x}\ _2 \leq \ \mathbf{x}\ _1 \leq \sqrt{d} \ \mathbf{x}\ _2</math>.</li> <li><math>\ A\mathbf{x}\  \leq \ A\ _2 \ \mathbf{x}\ </math>.</li> <li><math>\ A\ _2 \leq \ A\ _F</math>.</li> <li><b>Mean-value theorem</b> (<math>h</math> cont. on <math>[a, b]</math>, diff. on <math>(a, b)</math>):  <math display="block">h'(c) = \frac{h(b) - h(a)}{b - a}, \quad \exists c \in (a, b).</math> </li> <li><b>Fund. theorem of calculus</b> (<math>h</math> diff. on <math>[a, b]</math>, <math>h'</math> cont. on <math>[a, b]</math>):  <math display="block">h(b) - h(a) = \int_a^b h'(t) dt.</math> </li> </ul>	

- ( $L$ -smooth, convex,  $\gamma := \frac{1}{L}$ )  $f(\mathbf{x}_T) - f^* \leq \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^*\|^2$
- Sufficient decrease to bound gradients of VA with telescoping sum.
- ( $L$ -smooth,  $\mu$ -SC,  $\gamma := \frac{1}{L}$ )  $f(\mathbf{x}_T) - f^* \leq \frac{L}{2} \left(1 - \frac{\mu}{L}\right)^T \|\mathbf{x}_0 - \mathbf{x}^*\|^2$
- Use  $\mu$ -SC to strengthen VA bound for squared norm  $\Rightarrow$  Upper bound “noise” with  $f^* \leq f(\mathbf{x}_{t+1})$  and SD  $\Rightarrow$  Smoothness on  $(\mathbf{x}^*, \mathbf{x}_T)$ .

## Projected gradient descent

- **Update rule** ( $X \subset \mathbb{R}^d$  is closed and convex):  

$$\mathbf{y}_{t+1} = \mathbf{x}_t - \gamma_t \nabla f(\mathbf{x}_t)$$

$$\mathbf{x}_{t+1} = \Pi_X(\mathbf{y}_{t+1}) := \operatorname{argmin}_{\mathbf{x} \in X} \|\mathbf{x} - \mathbf{y}_{t+1}\|^2.$$
- ( $\mathbf{x} \in X, \mathbf{y} \in \mathbb{R}^d$ ):  $(\mathbf{x} - \Pi_X(\mathbf{y}))^\top (\mathbf{y} - \Pi_X(\mathbf{y})) \leq 0$ .
- Constrained 1st-order optimality  $\Rightarrow$  Rearrange.
- ( $\mathbf{x} \in X, \mathbf{y} \in \mathbb{R}^d$ ):  $\|\mathbf{x} - \Pi_X(\mathbf{y})\|^2 + \|\mathbf{y} - \Pi_X(\mathbf{y})\|^2 \leq \|\mathbf{x} - \mathbf{y}\|^2$ .
- Cosine theorem on previous.
- If  $\mathbf{x}_{t+1} = \mathbf{x}_t$ , then  $\mathbf{x}_t = \mathbf{x}^*$ .
- $(\mathbf{x} - \mathbf{x}_t)^\top (\mathbf{y}_{t+1} - \mathbf{x}_t) = (\mathbf{x} - \mathbf{x}_{t+1})^\top (\mathbf{y}_{t+1} - \mathbf{x}_{t+1}) \leq 0, \forall \mathbf{x} \in X$  to show that constrained 1st-order optimality holds.
- **Projected SD**:
- 
- ( $L$ -smooth, convex,  $\gamma := \frac{1}{L}$ ):
- 

## Coordinate descent

- **Coordinate-wise SD**:
- 
- **Convergence results** ( $\mu$ -PL,  $\mathcal{L}$ -CS,  $\bar{L} = \frac{1}{d} \sum_{i=1}^d L_i$ ):
  - ( $i \sim \text{Unif}([d])$ )
  - 
  - ( $i \sim \text{Cat}(L_1/\sum_{j=1}^d L_j, \dots, L_d/\sum_{j=1}^d L_j)$ )
  - 
  - ( $i \in \operatorname{argmax}_{j \in [d]} |\nabla_j f(\mathbf{x}_t)|$ )

## Nonconvex functions

- ( $L$ -smooth):
- 
- **Trajectory analysis**: Optimize  $f(\mathbf{x}) := \frac{1}{2} \left( \prod_{k=1}^d x_k - 1 \right)^2$ .
- $\frac{\partial f(\mathbf{x})}{\partial x_i} = (\prod_k x_k - 1) \prod_{k \neq i} x_k$  ( $\nabla f(\mathbf{x}) = \mathbf{0}$  if 2 dims are 0 or all 1).
- $\frac{\partial^2 f(\mathbf{x})}{\partial x_i^2} = \left( \prod_{k \neq i} x_k \right)^2$ .
- $\frac{\partial^2 f(\mathbf{x})}{\partial x_i \partial x_j} = 2 \prod_{k \neq i} x_k \prod_{k \neq j} x_k - \prod_{k \neq i, j} x_k$ , if  $i \neq j$ .
- $c$ -**balanced**: Let  $\mathbf{x} > \mathbf{0}$ ,  $c \geq 1$ .  $\mathbf{x}$  is  $c$ -balanced if  $x_i \leq c \cdot x_j, \forall i, j \in [d]$ .
- If  $\mathbf{x}_t$  is  $c$ -balanced,  $\gamma > 0$ , then  $\mathbf{x}_{t+1}$  is  $c$ -balanced and  $\mathbf{x}_{t+1} \geq \mathbf{x}_t$ .
- 
- If  $\mathbf{x}$  is  $c$ -balanced, then for any  $I \subseteq [d]$ , we have
$$\prod_{k \notin I} x_k \leq c^{|I|} \left( \prod_{k=1}^d x_k \right)^{1 - |I|/d} \leq c^{|I|}.$$
- 
- Let  $\mathbf{x}$  be  $c$ -balanced and  $\prod_k x_k \leq 1$ , then
$$\|\nabla^2 f(\mathbf{x})\|_2 \leq \|\nabla^2 f(\mathbf{x})\|_F \leq 3dc^2.$$
Thus,  $f$  is smooth along the whole trajectory of GD.
- 
- **Convergence** ( $\gamma = \frac{1}{3dc^2}$ ,  $\mathbf{x}_0 > \mathbf{0}$  and  $c$ -balanced,  $\delta \leq \prod_k x_{0,k} < 1$ )
$$f(\mathbf{x}_T) \leq \left(1 - \frac{\delta^2}{3e^4}\right)^T f(\mathbf{x}_0).$$

- 
- $\delta$  decays polynomially in  $d$ , so we must start  $\mathcal{O}(1/\sqrt{d})$  from  $\mathbf{x}^* = \mathbf{1}$ .

## Frank-Wolfe

- $\text{LMO}_X(\mathbf{g}) := \operatorname{argmin}_{\mathbf{z} \in X} \mathbf{g}^\top \mathbf{z}$ .
- **Update rule**:
$$\mathbf{s}_t = \text{LMO}_X(\nabla f(\mathbf{x}_t))$$

$$\mathbf{x}_{t+1} = (1 - \gamma_t) \mathbf{x}_t + \gamma_t \mathbf{s}_t.$$
- If  $X = \operatorname{conv}(\mathcal{A})$ , then  $\text{LMO}_X(\mathbf{g}) \in \mathcal{A}$ .
- Advantages: (1) Iterates are always feasible if  $X$  is convex, (2) No projections, (3) Iterates have simple sparse representations as convex combination of  $\{\mathbf{x}_0, \mathbf{s}_0, \dots, \mathbf{s}_t\}$ .
- LMO of unit  $\ell_1$ -ball:  $\text{LMO}(\mathbf{g}) = -\operatorname{sgn}(g_i) \mathbf{e}_i, i \in \operatorname{argmax}_{j \in [d]} |g_j|$ .
- **Optimality gap**:  $g(\mathbf{x}) := \nabla f(\mathbf{x})^\top (\mathbf{x} - \mathbf{s}), \mathbf{s} = \text{LMO}_X(\nabla f(\mathbf{x}))$ .
- $g(\mathbf{x}) \geq f(\mathbf{x}) - f^*$ .
- 
- **Descent lemma**:  $f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \gamma_t g(\mathbf{x}_t) + \gamma_t^2 \frac{L}{2} \|\mathbf{s}_t - \mathbf{x}_t\|^2$ .
- **Convergence** ( $L$ -smooth, convex,  $X$  is compact,  $\gamma_t = \frac{2}{t+2}$ ):
$$f(\mathbf{x}_T) - f^* \leq \frac{2L}{T+1} \operatorname{diam}(X)^2.$$
- Lemma  $-f^* \Rightarrow$  Use  $g(\mathbf{x}) \geq f(\mathbf{x}) - f^* \Rightarrow$  Rearrange and induction.
- **Affine equivalence**:  $(f, X)$  and  $(f', X')$  are affinely equivalent if  $f'(\mathbf{x}) = f(A\mathbf{x} + \mathbf{b})$  and  $X' = \{A^{-1}(\mathbf{x} - \mathbf{b}) \mid \mathbf{x} \in X\}$ . Then,
$$\nabla f'(\mathbf{x}') = A^\top \nabla f(\mathbf{x}), \quad \mathbf{x} = A^{-1}(\mathbf{x}' - \mathbf{b})$$

$$\text{LMO}_{X'}(\nabla f'(\mathbf{x}')) = A^{-1}(\mathbf{s} - \mathbf{b}), \quad \mathbf{s} = \text{LMO}_X(\nabla f(\mathbf{x})).$$

- **Curvature constant**:
$$C_{(f,X)} := \sup_{\substack{\mathbf{x}, \mathbf{s} \in X, \gamma \in (0,1] \\ \mathbf{y} = (1-\gamma)\mathbf{x} + \gamma\mathbf{s}}} \frac{1}{\gamma^2} (f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})).$$

- **Affine invariant convergence**:  $f(\mathbf{x}_T) - f^* \leq \frac{4C_{(f,X)}}{T+1}$ .

- Descent lemma w.r.t.  $C_{(f,X)}$  by setting  $\mathbf{x} = \mathbf{x}_t, \mathbf{s} = \text{LMO}_X(\nabla f(\mathbf{x}_t))$  in the supremum.

- **Convergence of  $g(\mathbf{x}_t)$** :  $\min_{1 \leq t \leq T} g(\mathbf{x}_t) \leq \frac{27/2 \cdot C_{(f,X)}}{T+1}$ .

## Newton's method

- **Update rule**:  $\mathbf{x}_{t+1} = \mathbf{x}_t - \nabla^2 f(\mathbf{x}_t)^{-1} \nabla f(\mathbf{x}_t)$ .
- **Interp. 1**: Adaptive gradient descent.
- **Interp. 2**: Minimizes second-order Taylor approximation around  $\mathbf{x}_t$ :
$$\mathbf{x}_{t+1} \in \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^\top (\mathbf{x} - \mathbf{x}_t) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_t)^\top \nabla^2 f(\mathbf{x}_t) (\mathbf{x} - \mathbf{x}_t).$$
- **Convergence** ( $\|\nabla^2 f(\mathbf{x})^{-1}\| \leq \frac{1}{\mu}, \|\nabla^2 f(\mathbf{x}) - \nabla^2 f(\mathbf{y})\| \leq B \|\mathbf{x} - \mathbf{y}\|$ ):
$$\|\mathbf{x}_{t+1} - \mathbf{x}^*\| \leq \frac{B}{2\mu} \|\mathbf{x}_t - \mathbf{x}^*\|^2.$$
- $\mathbf{x}_{t+1} - \mathbf{x}^* \leq \mathbf{x}_t - \mathbf{x}^* + H(\mathbf{x}_t)^{-1} (\nabla f(\mathbf{x}^*) - \nabla f(\mathbf{x}_t)) \Rightarrow h(t) := \nabla f(\mathbf{x} + t(\mathbf{x}^* - \mathbf{x}))$  with fundamental theorem of calculus  $\Rightarrow$  Take norm of both sides and simplify using  $\|A\mathbf{x}\| = \|A\|_2 \|\mathbf{x}\|$  and assumptions.
- Ensure bounded inverse Hessians by requiring strong convexity over  $X$ .
- If  $\|\mathbf{x}_0 - \mathbf{x}^*\| \leq \frac{\mu}{B}$ , then  $\|\mathbf{x}_T - \mathbf{x}^*\| \leq \frac{\mu}{B} \left(\frac{1}{2}\right)^{2^T - 1}$ .

## Quasi-Newton methods

### Subgradient method

### Mirror descent

### Stochastic optimization

### Variance reduction

### Min-max optimization

### Variational inequalities