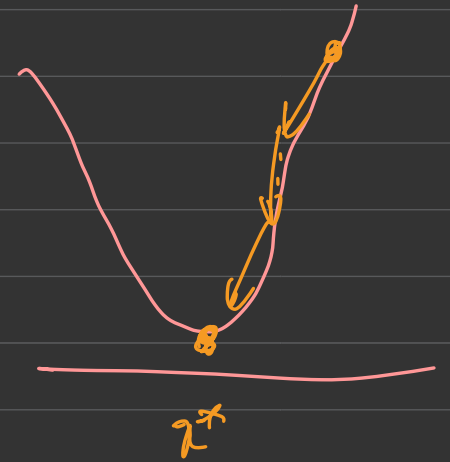


- Descent methods : Gradient, Steepest descent
- Gradient descent convergence :
 - Vanilla analysis (average iterates)
 - Convex & Lipschitz (average iterates $O(1/\sqrt{k})$)
 - Convex & L -smooth [last iterates $O(1/k)$]

Differentiable function minimization (unconstrained)

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\begin{array}{ll} \text{minimize} & f(\underline{x}) \\ \text{subject to} & \underline{x} \in \mathbb{R}^n \end{array}$$



Iterative "descent" algorithms:

Start with \underline{x}_0 and compute a sequence $\underline{x}_1, \underline{x}_2, \dots$

Such that

$$f(\underline{x}_{t+1}) < f(\underline{x}_t) \quad t=0, 1, \dots$$

Descent direction:

$$f(\underline{x}_t + \underline{d}_t) \approx f(\underline{x}_t) + \nabla f^T(\underline{x}_t) \underline{d}_t$$

directional derivative: Change in f for a small step \underline{d}_t

$$\begin{aligned} f'(\underline{x}_t; \underline{d}_t) &:= \lim_{\tau \rightarrow 0} \frac{f(\underline{x}_t + \tau \underline{d}_t) - f(\underline{x}_t)}{\tau} \\ &= \nabla f^T(\underline{x}_t) \underline{d}_t \end{aligned}$$

• \underline{d}_t is a descent direction at \underline{x}_t if $\nabla f^T(\underline{x}_t) \underline{d}_t < 0$

• For convex f , (first-order condition)

$$\nabla f^T(\underline{x}_t) (\underline{x}_{t+1} - \underline{x}_t) \geq 0 \Rightarrow f(\underline{x}_{t+1}) \geq f(\underline{x}_t)$$

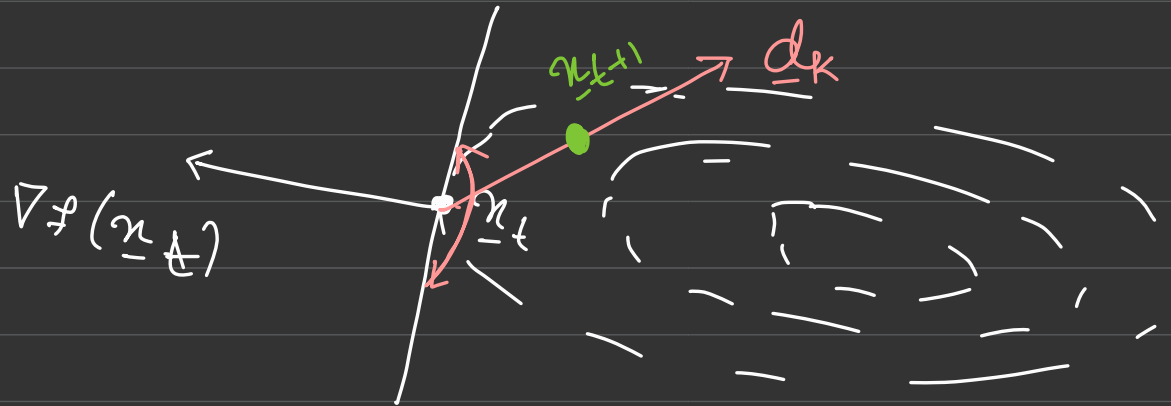
So a descent direction must satisfy

$$\nabla f^T(\underline{x}_t) \underline{d}_t \leq 0$$

Iterative descent methods:

Step size

$$\underline{x}_{t+1} = \underline{x}_t + \eta_t \underline{d}_t ; \eta_t > 0$$



Gradient descent: [Cauchy 1847]

$$\underline{d}_t = -\beta \nabla f(\underline{x}_t)$$

$\beta > 0$

$$\underline{x}_{t+1} = \underline{x}_t - \eta_t \nabla f(\underline{x}_t)$$

$$\underline{d}_t = -\nabla f(\underline{x}_t) \Rightarrow \nabla f^T(\underline{x}_t) (-\nabla f(\underline{x}_t)) = -\|\nabla f(\underline{x}_t)\|_2^2 < 0$$

Steepest descent method:

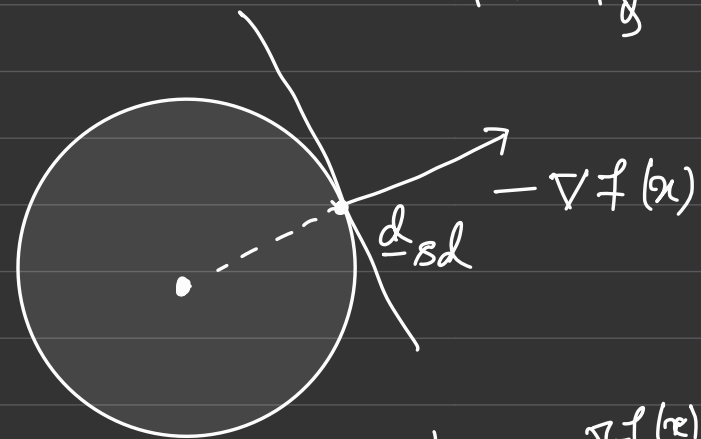
→ descent direction that is as negative as possible that yields the greatest rate of objective value improvement

$$\underline{d}_{sd} = \arg \min_{\underline{d}} \{ \nabla f^T(\underline{x}) \underline{d} : \|\underline{d}\| \leq 1 \}$$

$\|\underline{d}\|_2 \leq 1$

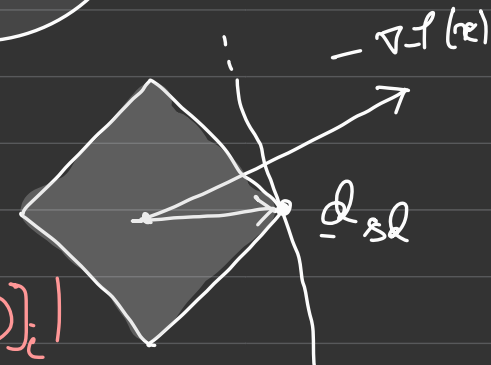
For Euclidean norm:

$$\underline{d}_{sd} = -\nabla f(\underline{x})$$



For l_1 -norm:

$$\underline{d}_{sd} = \arg \min_{\underline{d}} \{ \nabla f^T(\underline{x}) \underline{d} : \|\underline{d}\|_1 \leq 1 \}$$



Let i be any index s.t. $\|\nabla f(\underline{x})\|_\infty = |[\nabla f(\underline{x})]_i|$

Then $\underline{d}_{sd} = -\text{sign}\left(\frac{\partial f(\underline{x})}{\partial x_i}\right) \underline{e}_i$: coordinate descent.

Vanilla analysis:

$$\underline{x}_{t+1} = \underline{x}_t - \eta \nabla f(\underline{x}_t)$$

$$t = 0, \dots, T-1$$

Define : $\underline{g}_t = \nabla f(\underline{x}_t)$

so $\underline{g}_t = (\underline{x}_{t+1} - \underline{x}_t) / \eta$

Let us relate this vector to our current direction from an optimum \underline{x}^* : $\underline{x}_t - \underline{x}^*$

$$\underline{g}_t^T (\underline{x}_t - \underline{x}^*) = \frac{1}{\eta} (\underline{x}_{t+1} - \underline{x}_t)^T (\underline{x}_t - \underline{x}^*)$$

Cosine theorem: $2 \underline{v}^T \underline{w} = \|\underline{v}\|^2 + \|\underline{w}\|^2 - \|\underline{v} - \underline{w}\|^2$

$$\underline{g}_t^T (\underline{x}_t - \underline{x}^*) = \frac{1}{2\eta} \left[\|\underline{x}_{t+1} - \underline{x}_t\|_2^2 + \|\underline{x}_t - \underline{x}^*\|_2^2 - \|\underline{x}_{t+1} - \underline{x}^*\|_2^2 \right]$$

⊗

$$= \frac{1}{2\eta} \left[\eta^2 \|\underline{g}_t\|_2^2 + \|\underline{x}_t - \underline{x}^*\|_2^2 - \|\underline{x}_{t+1} - \underline{x}^*\|_2^2 \right]$$

$$= \frac{\eta}{2} \|\underline{g}_t\|_2^2 + \frac{1}{2\eta} \left[\|\underline{x}_t - \underline{x}^*\|_2^2 - \|\underline{x}_{t+1} - \underline{x}^*\|_2^2 \right]$$

⊗ Sum over the iteration t .

$$\boxed{a_n - a_{n+1}}$$

$$\sum_{t=0}^{T-1} \underline{g}_t^T (\underline{x}_t - \underline{x}^*) = \frac{\eta}{2} \sum_{t=0}^{T-1} \|\underline{g}_t\|^2$$

$$+ \frac{1}{2\eta} \left[\|\underline{x}_0 - \underline{x}^*\|_2^2 - \|\underline{x}_T - \underline{x}^*\|_2^2 \right]$$

$$\leq \frac{\eta}{2} \sum_{t=0}^{T-1} \|\underline{g}_t\|^2 + \frac{1}{2\eta} \|\underline{x}_0 - \underline{x}^*\|_2^2$$

Now, suppose f is convex: $f(\underline{y}) \geq f(\underline{x}) + \nabla f^T(\underline{x})(\underline{y} - \underline{x})$

$$\left. \begin{array}{l} \underline{y} = \underline{x}^* \\ \underline{x} = \underline{x}_t \end{array} \right\} \Rightarrow f(\underline{x}_t) - f(\underline{x}^*) \leq \nabla f^T(\underline{x}_t) (\underline{x}_t - \underline{x}^*) = \underline{g}_t^T (\underline{x}_t - \underline{x}^*)$$

• Upper bound on the average error: $f(\underline{x}_t) - f(\underline{x}^*)$

$$\sum_{t=0}^{T-1} f(\underline{x}_t) - f(\underline{x}^*) \leq \frac{\eta}{2} \sum_{t=0}^{T-1} \|\underline{g}_t\|^2 + \frac{1}{2\eta} \|\underline{x}_0 - \underline{x}^*\|_2^2$$

• Last iterate is not necessarily the best one as "fixed" step size can make steps overshoot & increase function value

• For Lipschitz convex functions:

$f: X \rightarrow \mathbb{R}^n$ is called Lipschitz continuous if there exists $B \geq 0$

$$\|f(x) - f(y)\| \leq B \|x - y\| \quad \forall x, y \in X$$

$$\Leftrightarrow \|\nabla f(x)\| \leq B \quad \forall x \in X$$

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\|x_0 - x^*\| \leq R$$

Then,

$$\sum_{t=0}^{T-1} (f(x_t) - f(x^*)) \leq \frac{\eta}{2} B^2 T + \frac{1}{2\eta} R^2 \quad (**)$$

So, choose η such that $g(\eta) = \frac{\eta}{2} B^2 T + \frac{1}{2\eta} R^2$ is minimized.

$$\frac{d}{d\eta} g(\eta) = 0 \Rightarrow \frac{1}{2} B^2 T - \frac{1}{2\eta^2} R^2 = 0$$

$$\Rightarrow \eta = \frac{R}{B\sqrt{T}} \quad \text{and } g(\eta) = RB\sqrt{T}$$

$$(**) \div T$$

$$\frac{1}{T} \sum_{t=0}^{T-1} (f(x_t) - f(x^*)) \leq \frac{RB}{\sqrt{T}} \approx O\left(\frac{1}{\sqrt{T}}\right)$$

independent of n

but $R \neq B$

does depend on η

So to obtain $\min_{t=0..T-1} f(x_t) - f(x^*) \leq \epsilon$

we need $T \geq \frac{R^2 B^2}{\epsilon^2}$. E.g. $\epsilon = 10^{-8}$
 R and $B \approx 10^8$??