

Dual norm

Let  $\|\cdot\|$  be any norm on  $\mathbb{R}^n$

$\|\cdot\|_*$  is a dual norm of  $\|\cdot\|$ , defined on  $\mathbb{R}^n$  as

follows:

$$\forall z \in \mathbb{R}^n, \quad \|z\|_* = \sup \{ z^T v : \|v\| \leq 1 \}$$

How much can the vector  $z$  inflate  $v$  when  $\|v\| \leq 1$ ?

$$z^T v \leq \|z\|_* \|v\|$$

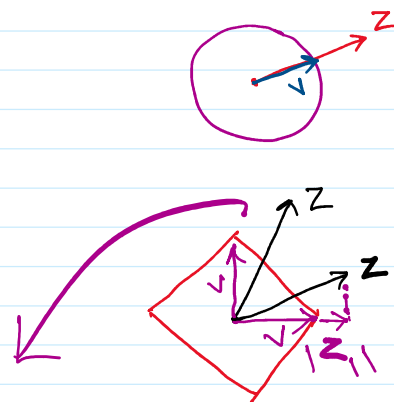
E.g.,  $\|\cdot\| = \|\cdot\|_2$  ( $\ell_2$  norm)

$$\sup (z^T v) = \|z\|_2$$

Say,  $\|\cdot\| = \|\cdot\|_1$

$$z^T v = |z_1|$$

$$z^T v = |z_2|$$



$$\|\cdot\|_* \text{ when } \|\cdot\| = \|\cdot\|_1 \Rightarrow \|\cdot\|_* = \|\cdot\|_\infty$$

Dual norms of  $\|\cdot\|_p$  when:

$$p = 2 \Rightarrow \|\cdot\|_* = \|\cdot\|_2$$

$$p = 1 \Rightarrow \|\cdot\|_* = \|\cdot\|_\infty$$

$$p = \infty \Rightarrow \|\cdot\|_p = \|\cdot\|_1$$

Dual norm of  $\ell_p$  is  $\ell_q$  norm

$$\text{where } \frac{1}{p} + \frac{1}{q} = 1 \Leftrightarrow q = \frac{p}{p-1}$$

$$(\|\cdot\|_*)_* = \|\cdot\|$$

## Steepest Descent

Assume  $f(x)$  attains its minimum

$$x^* \in \arg \min_{x \in \text{dom} f} f(x)$$

$$\text{Descent method: } x^{(k+1)} = x^{(k)} + t^{(k)} \Delta x^{(k)}$$

where  $\Delta x^{(k)}$  is a descent direction

$$\text{i.e. } -\nabla f(x^{(k)})^T \Delta x^{(k)} > 0$$

Remember: First-order approx. of  $f$  around  $x^{(k)}$

$$f(x^{(k+1)}) \approx f(x^{(k)}) + t^{(k)} \nabla f(x^{(k)})^T \Delta x^{(k)}$$

when  $\Delta x^{(k)}$  is a descent direction

$$\Rightarrow f(x^{(k+1)}) < f(x^{(k)})$$

The reduction in function value  $\propto \nabla f(x^{(k)})^T \Delta x^{(k)}$

$$\Rightarrow f(x^{(k+1)}) < f(x^{(k)})$$

The reduction in function value  $\propto \nabla f(x^{(k)})^T \Delta x^{(k)}$

Steepest descent is the descent method in which we have the most reduction in the objective function value (i.e.,  $\nabla f(x)^T \Delta x$  is the smallest)

↳ i.e.,  $-\nabla f(x)^T \Delta x$  is the largest value.

This is a non-rigorous statement.

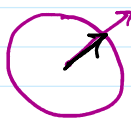
$$\begin{aligned} \text{we could ask } \Delta x_{sd} &= \arg \min_{v \in \mathbb{R}^n} \{ \nabla f(x)^T v \} \\ &= \arg \max_{v \in \mathbb{R}^n} \{ -\nabla f(x)^T v \} \end{aligned}$$

But this would always return  $\Delta x_{sd} \rightarrow -\begin{bmatrix} \infty \\ \vdots \\ \infty \end{bmatrix}$

Better idea: Normalized steepest descent direction:

$$\Delta x_{nsd} = \arg \min_{v \in \mathbb{R}^n} \{ \nabla f(x)^T v : \|v\| \leq 1 \}$$

$$\|\cdot\| = \|\cdot\|_2 \quad -\nabla f(x) = \arg \left( \max_{v \in \mathbb{R}^n} \{ -\nabla f(x)^T v : \|v\| \leq 1 \} \right)$$



$$\Delta x_{nsd} = \frac{-\nabla f(x)}{\|\nabla f(x)\|}$$

$$\begin{aligned} \text{Dual norm of } -\nabla f(x) &= \|\nabla f(x)\|_* \\ &= \|\nabla f(x)\|_2 \end{aligned}$$

$$\Delta x_{\text{nsd}} = \frac{-\nabla f(x)}{\|\nabla f(x)\|_2} = \|\nabla f(x)\|_*$$

$$\Delta x_{\text{sd}} = \|\nabla f(x)\|_* \Delta x_{\text{nsd}} \Rightarrow \text{Steepest descent direction}$$

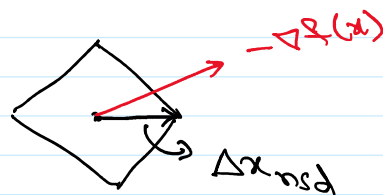
Steepest descent iteration:  $x^{(k+1)} = x^{(k)} + t^{(k)} \underbrace{\|\nabla f(x^{(k)})\|_*}_{\text{Steepest descent direction}} \Delta x_{\text{nsd}}^{(k)}$

When  $\|\cdot\| = \|\cdot\|_2 \Rightarrow \|\cdot\|_* = \|\cdot\|_2$   ~~$\|\nabla f(x^{(k)})\|_2 \cdot \frac{-\nabla f(x^{(k)})}{\|\nabla f(x^{(k)})\|_2}$~~

$$\Delta x_{\text{nsd}} = \frac{-\nabla f(x)}{\|\nabla f(x)\|_2} \Rightarrow$$

$\Rightarrow \text{SD} = \text{GD}$  when we are using  $\|\cdot\|_2$ .

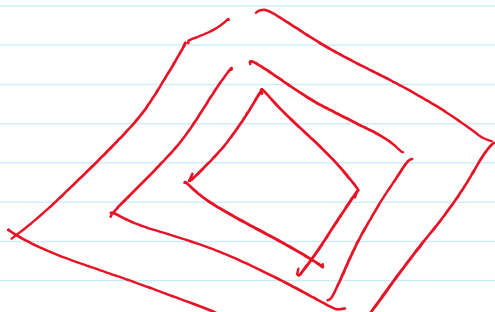
E.g.  $\|\cdot\| = \|\cdot\|_1$



In that case, the direction is coordinate in which  $-\frac{\partial f}{\partial x_i}$  is the largest.

The norm that one should use in steepest descent should be such that norm-ball has geometry that matches the geometry of the level sets.

Say:



when the norm is correct: SD converges faster than GD.  
 when the guess is wrong: SD converges slower than GD.

Special Case of Steepest Descent: Use  $P$ -quadratic norm, where  $P \succ 0$ .  $P \in \mathbb{R}^{n \times n}$

Reminder:  $\|v\|_P = (v^T P v)^{1/2} = \|P^{1/2} v\|_2$

Also called Gradient descent with Preconditioning.  
 $P$  is called the preconditioning matrix.

Facts:  $\|z\|_*$  when  $\|\cdot\| = \|\cdot\|_P$

$$\|z\|_* = \|P^{-1/2} z\|_2 = (z^T P^{-1} z)^{1/2}$$

$$\Delta x_{\text{nsd}} = \frac{-P^{-1} \nabla f(x)}{\|\nabla f(x)\|_*}$$

$$= -(\nabla f(x)^T P^{-1} \nabla f(x))^{-1/2} P^{-1} \nabla f(x)$$

$$\Delta x_{\text{sd}} = \|\nabla f(x)\|_* \Delta x_{\text{nsd}} = -P^{-1} \nabla f(x)$$

Bottom line: SD with  $\|\cdot\|_P$  is

$$x^{(k+1)} = x^{(k)} - t^{(k)} P^{-1} \nabla f(x) \quad \text{---} \quad (*)$$

$$GD \Rightarrow P = I ;$$

$$\text{Newton's method} \Rightarrow P = \nabla^2 f(x)$$

Steepest descent under  $P$ -quadratic norm is equivalent to doing a change of coordinates of our problem by  $P^{-1/2}$  and then running gradient descent.

$$\boxed{f(x)} \text{ with } x \in \mathbb{R}^n$$

$$\text{Let } \bar{x} = P^{1/2} x ; \quad \bar{f}(\bar{x}) = \boxed{f(P^{-1/2} x)}$$

$$\Leftrightarrow x = P^{-1/2} \bar{x}$$



$$\nabla_{\bar{x}} \bar{f}(\bar{x}) = P^{-1/2} \nabla f(P^{-1/2} x)$$

$$x^{(u+1)} = x^{(u)} + t^{(u)} \underbrace{\nabla f(x^{(u)})^T}_{\nabla f^T P^{-1/2}} \underbrace{\Delta x^{(u)}}_{\nabla f(x)} \Delta x^{(u)}$$

$$\nabla f^T P^{-1/2} \nabla f(x) = P^{-1/2} \nabla f$$

$$\nabla f^T P^{-1} \nabla f = \|\nabla f(x)\|_x^2$$

Special Case:  $f(x) = x^T Q x \Rightarrow \text{cond}(f(x)) = \frac{\lambda_{\max}(Q)}{\lambda_{\min}(Q)}$

Condition number of  $f(x)$  is determined by eigenvalues of  $Q$ .

$$PICK = P = Q$$

$$\Rightarrow P^{-1/2} = Q^{-1/2}$$

$$\tilde{P}(\tilde{x}) = \tilde{P}(P^{-1/2}x) = (P^{-1/2}x)^T Q (P^{-1/2}x)$$

$$= x^T \underbrace{P^{-1/2} Q P^{-1/2}}_{= I} x$$

$$= I$$

$$= x^T x$$

$$\text{cond}(\tilde{P}(\tilde{x})) = 1$$