

What is optimization?

consider a function $f: \underbrace{\mathbb{R}^n}_{\text{dom } f} \rightarrow \mathbb{R}$

$$\text{e.g., } f(x) = \log(x)$$

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

$$\text{dom } f = (0, \infty)$$

optimization is concerned with finding the maximum or minimum value that f takes:

$$\min_{x \in \mathbb{R}^n} f(x) \quad \text{or} \quad \max_{x \in \mathbb{R}^n} f(x)$$

↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓

If min or max exists infimum supremum

$$\text{e.g., } f(x) = \log(x)$$

$$\min_x f(x) \Rightarrow \text{does not exist}$$

$$f(x) = x^2 \Rightarrow \min_x f(x) = 0$$

$$x_n = \frac{1}{n} : n=1, 2, \dots$$

$$\min_n x_n \text{ does not exist}$$

$$\text{but } \inf_n x_n = 0$$

(or) Finding the x that gives the minimum or maximum value of f

$$\boxed{\arg \min_{x \in \mathbb{R}^n} f(x)} \quad \text{or}$$

$$\boxed{\arg \max_{x \in \mathbb{R}^n} f(x)}$$

$$\underset{x \in \mathbb{R}^n}{\text{argmin}} f(x) \quad \text{or}$$

$$\underset{x \in \mathbb{R}^n}{\text{argmax}} f(x)$$

e.g.: $f(x) = x^2$

$$\min_x f(x) = 0$$

$$\underset{x}{\text{argmin}} f(x) = 0$$

$$f(x) = x^2 + 1$$

$$\min_x f(x) = 1$$

$$\underset{x}{\text{argmin}} f(x) = 0$$

$$f(x) = (x-1)^2$$

$$\min_x f(x) = 0$$

$$\underset{x}{\text{argmin}} f(x) = 1$$

Practically, the function $f(x)$ represents some real-world quantity that we care about and $x \in \mathbb{R}^n$ corresponds to n parameters of the problem that alter the quantity of interest.

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

Example: Portfolio optimization

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

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

→ investment in stock 1
→ investment in stock n

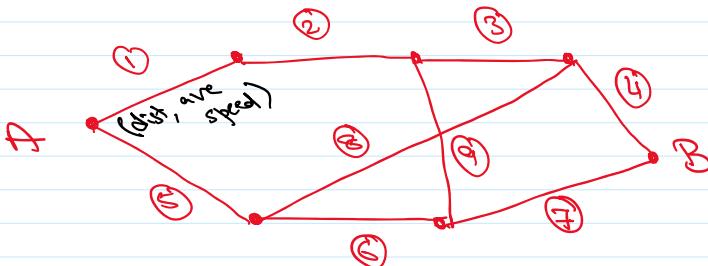
\hookrightarrow Profit / loss after 6 months

$$\arg \max_x f(x)$$

Constraints: Total budget $\Rightarrow \sum_i x_i \leq B$

$$x_i \geq 0$$

Example: Route optimization w/ respect to minimization of time.



option 1:

$$① \rightarrow ② \rightarrow ③ \rightarrow ④$$

option 2

$$① \rightarrow ② \rightarrow ⑤ \rightarrow ⑦$$

$f(x) =$ Time needed to go from A to B

\hookrightarrow encoding the path

$$f(x): \{0, 1\}^9 \rightarrow \mathbb{R}$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_9 \end{bmatrix} \Rightarrow \begin{array}{ll} x_i = 0 & \text{if } i^{\text{th}} \text{ edge is not taken} \\ & \\ x_i = 1 & \text{if } i^{\text{th}} \text{ edge is taken} \end{array}$$

$$x = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix} \Rightarrow \text{option 1}$$

$$f(x) = \frac{\text{dist}(1)}{\text{ave. Speed}(1)} \times x_1 + \frac{\text{dist}(2)}{\text{ave. Speed}(2)} \times x_2 + \dots + \frac{\text{dist}(9)}{\text{ave. Speed}(9)} \times x_9$$

\hookrightarrow time taken to traverse edge 1

This problem involves optimizing a 9-dimensional binary vector $x \in \{0,1\}^9$

$$\underset{x \in \{0,1\}^9}{\operatorname{argmin}} f(x)$$

and

$$\min_x f(x)$$

↳ This actually is not over all $\{0,1\}^9$

We have constraints that x must correspond to a valid route from A to B

$$X = \left\{ x : x \text{ corresponds to a valid route from A to B.} \right\}$$

Much Smaller than 2^9

$$\underset{x \in X}{\operatorname{argmin}} f(x)$$

↳ Constraint set in optimization

↳ This is an example of combinatorial optimization.

Other Examples

① IC design layout and routing of wires

② Machine learning

③ Robot motion planning

Summary

Optimization problems are everywhere

- ① Modeling the problem into an optimization framework
that is solvable.
→ Theory
- ② Recognizing the nature of the optimization framework
and manipulating it so that it is efficient to solve.
→ Theory / Implementation
- ③ Using an appropriate numerical optimization method
to help us solve the problem ⇒ Run-time
→ Theory
- ④ Understanding the 'goodness' of the solution provided by the algorithm.



Mathematical optimization: Fundamental concepts

- ① Unless otherwise stated, $f(\cdot)$ is assumed to attain its minimum or maximum value.
- ② Without loss of generality, we will stick with $\min_x f(x)$ or $\max_x f(x)$, rather than $\arg \min_x f(x)$ or $\arg \max_x f(x)$

An optimization problem will be written as

$$\min_{x \in \mathbb{R}^n} f(x)$$

$$\text{or } \max_{x \in \mathbb{R}^n} f(x)$$

but when the variable is obvious then we write

$$\min f(x) \quad \text{or} \quad \max f(x)$$

$f \Rightarrow$ Objective function

$x \Rightarrow$ Optimization Variable

Unconstrained optimization

$$\min_{x \in \mathbb{R}^n} f(x) \quad \text{or} \quad \min_{x \in \text{dom}} f(x)$$

$\hookrightarrow 'x'$ should be searched over the entire domain of f .

Solution \Rightarrow A solution to an optimization problem

$\min_{x \in \mathbb{R}^n} f(x)$ is often denoted by $x^* \in \mathbb{R}^n$ and it

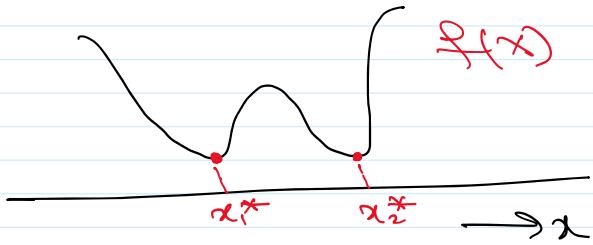
satisfies :

$$f(x^*) \leq f(x) \quad \forall x \in \mathbb{R}^n$$

x^* is referred to as a minimizer or

Optimal value.

We can have multiple solutions ; x_1^*, x_2^*, \dots



We can have 'no' solution.

e.g., $f(x) : (0, \infty) \rightarrow \mathbb{R}$

$$f(x) = -\frac{1}{x}$$

$\min_x f(x) \Rightarrow$ no x is going to give us the minimum value.

multiple Solutions

$$\Rightarrow f(x) = \cos(x)$$

$$\min_x \cos(x)$$

$$x^* \in \{\pm\pi, \pm 3\pi, \pm 5\pi, \dots\}$$

* An optimization algorithm is a numerical method that helps us numerically find a value of x close to an optimal x^* in a reasonable amount of time as a function of 'n' and other problem parameters.

↳ Mathematical programming

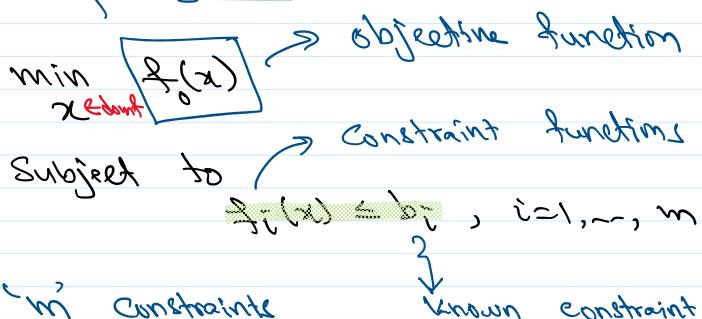
* All optimization problems that are min or max can be thought of as minimization problems only.

$$\max_{x \in \mathbb{R}^n} f(x) \iff \min_{x \in \mathbb{R}^n} -f(x)$$

$$\arg \max_x f(x) = \arg \min_x -f(x) = \min_x$$

$\min_x \max_y f(x, y) \Rightarrow$ not the focus of this course

Constrained optimization



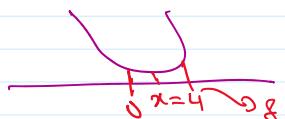
$\iff \min_{x \in \mathcal{X}} f_0(x) : \mathcal{X}$ is a subset of dom

$$\mathcal{X} = \{x : f_i(x) \leq b_i, i=1, \dots, m\}$$

e.g. $\min_{x \in \mathcal{X}} (x-4)^2$

E.g.

$$\min_{0 \leq x \leq 8} (x-4)^2$$



$$\mathcal{X} = \{x : x \in [0, 8]\}$$

$$\begin{aligned} 0 &\leq x \leq 8 \\ \Rightarrow x &\leq 8 \\ \Rightarrow -x &\leq 0 \end{aligned}$$

$$f_1(x) = x \Rightarrow b_1 = 8$$

$$f_2(x) = -x \Rightarrow b_2 = 0$$

$$\min_{x \in \text{dom } f} (x-4)^2$$

s.t.

$$\begin{aligned} x &\leq 8 \\ -x &\leq 0 \end{aligned}$$

Any $x \in \mathcal{X}$ is termed a 'feasible solution'.

In unconstrained optimization, all x 's are feasible.

Classes of optimization Problems

and the domain
of $f_0(x)$

Depending on $f_0(x)$, $f_i(x)$, $i=1, \dots, m$, we characterize an optimization problem to be from a certain class of problems, all of which can be solved using a particular set of numerical tools.

④ Combinatorial Optimization

$\Rightarrow x$ takes on values from a discrete set
(i.e., $\text{dom } f$ corresponds to a discrete set).

$$x \in \{0, 1\}^n$$

② Linear programming

The objective function $f_0(x)$ as well as the constraint functions $f_i(x), i=1, \dots, m$, are linear function

of x .



Integer programming

↪ Linear programs in which x takes only integer values.

③ Convex programming

↪ The function $f_0(x)$, the constraint functions $f_i(x) i=1, \dots, m$, are all convex.

A function f is convex if and only if

$$f(\alpha x + \beta y) \leq \alpha f(x) + \beta f(y)$$

* $\alpha, \beta \geq 0$, $\alpha + \beta = 1 \Rightarrow \alpha = 1 - \beta$

$$\Leftrightarrow f(\theta x + (1-\theta)y) \leq \theta f(x) + (1-\theta)f(y)$$

* $\theta \in [0,1]$

$x, y \in \text{dom } f$ and $\theta x + (1-\theta)y \in \text{dom } f$

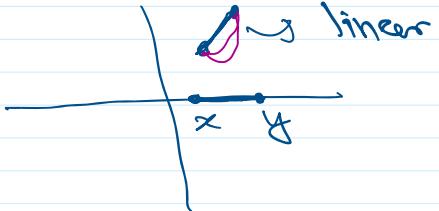
$\Leftrightarrow \text{dom } f$ has to satisfy the convexity property
that * $x, y \in \text{dom } f$

$$\theta x + (1-\theta)y \in \text{dom } f \quad * \quad \theta \in [0,1]$$

In other words, don't have to be a convex set.

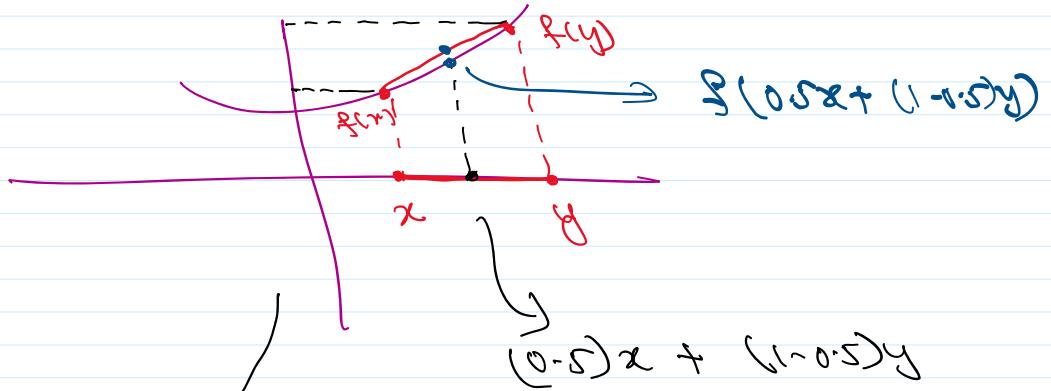
Linear function : $f(\alpha x + \beta y) = \alpha f(x) + \beta f(y)$

All linear functions are convex

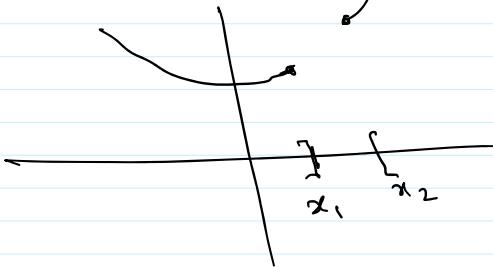


$$f(\theta x + (1-\theta)y) \leq \underline{\theta f(x) + (1-\theta)f(y)}$$

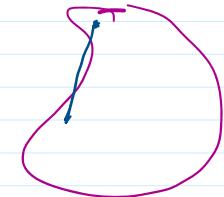
$$\theta \in [0,1]$$



E.g.,



convex set



not
a convex
set

④ Non Convex optimization

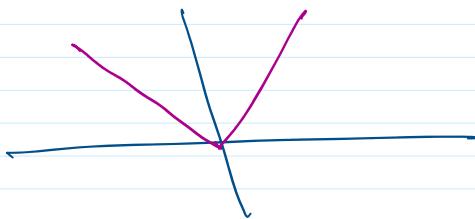
↳ Any optimization that is not convex

⑤ Nonlinear programming

↳ Any optimization that is not linear.

⑥ Nonsmooth optimization

↳ when $f_0(x)$ has discontinuous derivatives



⑦ Stochastic optimization

↳ when $f_0(x)$ or $f_i(x)$ have randomness involved in them.

Global optimization method

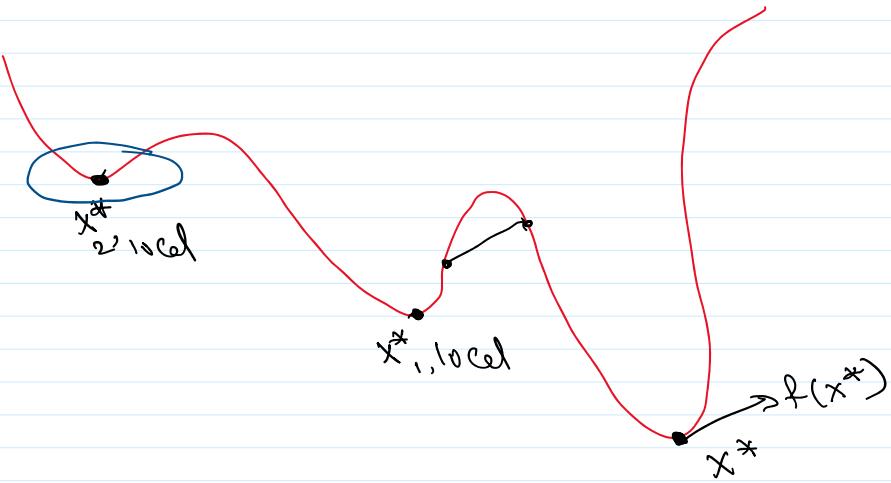
↳ A solver that provides an optimal solution x^* .

Convex optimization \Rightarrow we will see that global optimization is straight forward,

Nonconvex optimization \Rightarrow this is not computationally feasible in most cases (NP-hard).

Local optimization

↳ Find a solution x_{local}^* that is optimal in some neighborhood of x_{local}^* .



Why Convex Optimization?

- ① Many real-world problems tend to be convex.
- ② Convex optimization leads to global solutions.
- ③ The solvers are "efficient".
- ④ Many nonconvex problems can be 'relaxed' to convex problems and then solved globally.
- ⑤ Convex optim. can be used as an 'initialization' scheme for many nonconvex problems.

Class of Solvers

- ① Zeroth-order methods
 - ↳ solve the problem by having access to function values only.
- ② First-order methods
 - ↳ make use of first-order partial derivatives
 - ↳ ...

↪ Make use of first-order partial derivatives
of $f(x)$,

$$\frac{\partial f(x)}{\partial x_i}, i=1, \dots, n$$

↪ Use gradient information of the function

② Second-order methods

↪ Use both gradient information and second-order partial derivatives.

$$\frac{\partial^2 f(x)}{\partial x_i \partial x_j}, i, j = 1, \dots, n$$

↪ Use Hessian information

Examples of convex optimization

① Least Squares problem

$$f(x) = \|Ax - b\|_2^2 = \sum_{i=1}^n (a_i^T x - b_i)^2$$

$A: \mathbb{R}^{K \times n}$ matrix
 $b: \mathbb{R}^{n \times 1}$ vector

↪ Quadratic program

$b \Rightarrow$ observations vector (K observations)

$A \Rightarrow$ known matrix

Find closest x s.t.

$$\boxed{Ax \approx b}$$

If $K \geq n$ and A is full rank

$$x^* = (A^T A)^{-1} A^T b$$

↳ Analytical solution

↳ we can use convex optimization tools to solve this in $O(n^2 K)$ time.

Variants of LSWeighted LS

$$f(x) = \sum_{i=1}^k w_i (a_i^\top x - b_i)^2 \quad w_i \geq 0$$

$$= \|W(Ax - b)\|_2^2, \text{ where}$$

$$W = \text{diag}(w_1, w_2, \dots, w_k)$$

$$= \left\| \underbrace{WAx}_A - \underbrace{Wb}_b \right\|_2^2$$

Regularized LS

$$f(x) = \|Ax - b\|_2^2 + \rho \|x\|_2^2$$

↓
regularizer

regularization parameter
 $\rho \geq 0$

Ridge regression

Linear Programming

$$\min_x c^\top x \quad f_0(x) \text{ is linear}$$

$$\text{Subject to } \underbrace{a_i^\top x \leq b_i}_{\text{linear}} \quad i=1, \dots, m$$

$$f_1(x) \rightarrow f_m(x) \text{ linear}$$

Typical

$f_1(x) \rightarrow f_m(x)$ linear

$\mathcal{O}(n^2m)$ complexity if $m \geq n$

$$\min_x \max_{i=1, \dots, k} |a_i^\top x - b_i| = \min_x \|Ax - b\|_\infty$$

Chebyshev approximation

Compare to

$$\min_x \|Ax - b\|_2^2$$

Review of key linear algebra concepts

Inner product

Given $x, y \in \mathbb{R}^n$, the standard inner product is

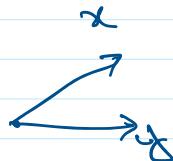
$$\langle x, y \rangle = x^\top y = \sum_{i=1}^n x_i y_i$$

and $x^\top x = \|x\|_2^2$, where $\|\cdot\|_2$ is Euclidean norm

Euclidean norm

Angle between two vectors: θ

★ $x^\top y = \|x\|_2 \|y\|_2 \cos(\theta)$



If $x^\top y = 0 \Leftrightarrow x \perp y$

If $x^\top y > 0 \Leftrightarrow x$ and y make an acute angle

Euclidean norm: $\|x\|_2 = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$

Cauchy-Schwarz Inequality

$$|x^T y| \leq \|x\|_2 \|y\|_2$$

Inner Product between matrices

Given $X, Y \in \mathbb{R}^{m \times n}$, inner product on matrices

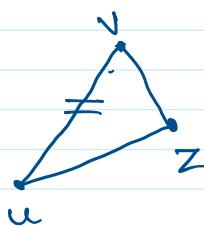
$$\langle X, Y \rangle = \text{trace}(X^T Y) = \sum_{i=1}^m \sum_{j=1}^n X_{ij} Y_{ij}$$

$$\langle X, X \rangle = \|X\|_F^2 = \sum_{i=1}^m \sum_{j=1}^n X_{ij}^2$$

What is a norm? focus on \mathbb{R}^n

A function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ with $\text{dom } f = \mathbb{R}^n$ is called a norm if:

- $f(x) \geq 0$ & $x \in \mathbb{R}^n$ — non-negativity
- $f(x) = 0$ only if $x = 0$ — definiteness
- $f(tx) = |t| f(x)$ & $t \in \mathbb{R}$, — homogeneity
- $f(x+y) \leq f(x) + f(y)$, — Triangle inequality



$$\|x\|_0$$

$\|x\|_2 \Rightarrow$ Euclidean norm $\Rightarrow l_2$ -norm

l_p -norm, $p \geq 1$

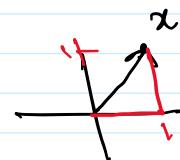
$$\|x\|_p = \left(|x_1|^p + |x_2|^p + \dots + |x_n|^p \right)^{\frac{1}{p}}$$

$$l_1\text{-norm} \Rightarrow \|x\|_1 = \sum_{i=1}^n |x_i|$$

$$p \rightarrow \infty$$

$$l_\infty\text{-norm} \Rightarrow \|x\|_\infty = \max_{i=1, \dots, n} |x_i|$$

↳ Chebyshev norm



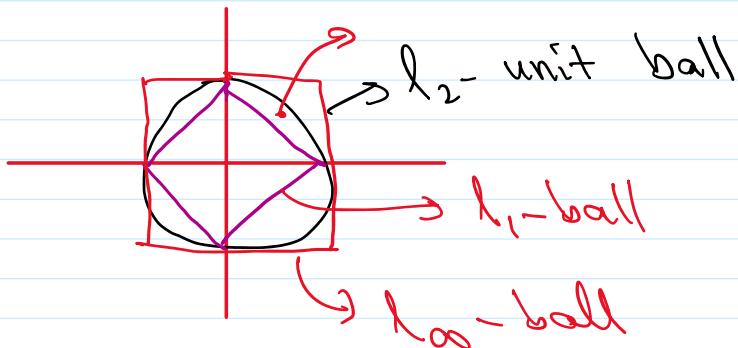
$$\|x\|_1 = 2$$

$$\|x\|_2 = \sqrt{2}$$

$$\|x\|_\infty = 1$$

Unit ball in R^n

$$B = \{x \in R^n, \|x\|_p \leq 1\}$$



All norms induce distances (metrics) on the vector space.

- $\text{dist}(x, y) = \|x - y\|_p$
- $\text{dist}(x, y) = \text{dist}(y, x)$
 - $\text{dist}(x, y) \geq 0 \quad \forall x, y$
 - $\text{dist}(x, y) = 0 \text{ only if } x = y$
 - $\text{dist}(x, y) \leq \text{dist}(x, z) + \text{dist}(z, y)$

Space of Symmetric matrices : $S^n \Rightarrow n \times n$ Symmetric matrices

Eigenvalue Decomposition

A matrix $M \in \mathbb{R}^{n \times n}$ said to have an eigenvalue decomposition if

$$M = \underline{U \Lambda U^{-1}} \Rightarrow M^T = \overbrace{\overbrace{U^T \Lambda^T U^T}^= U^T \Lambda U^T}$$

where $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$

$\lambda_1, \dots, \lambda_n \Rightarrow$ eigenvalues of M

Columns of $U \Rightarrow$ eigenvectors of $M \Rightarrow u_1, \dots, u_n$

If u_i is an eigenvector \Leftrightarrow

$$M u_i = \lambda_i u_i$$

Let $M \in S^n \Rightarrow$ Then M always has an EVD and the eigenvalues are always rel.

$$M^T = M \Leftrightarrow \underline{U \Lambda U^T} = M$$

T

T

$$v_1 = \dots \Leftrightarrow v_{n+1} = \dots$$

$$U^T U = U U^T = I$$

The eigenvectors are orthogonal to each other.

They are basis of \mathbb{R}^n .

$S_+^\gamma \subset S^\gamma$: S_+^γ : Positive Semidefinite $\xrightarrow{\text{symmetric}}$ matrices
 $\hookrightarrow \lambda_1, \dots, \lambda_n \geq 0 \xrightarrow{\text{PSD}}$

$S_{++}^\gamma \subset S^\gamma$: S_{++}^γ : Positive definite $\xrightarrow{\text{symmetric}}$ matrices
 $\lambda_1, \lambda_2, \dots, \lambda_n > 0 \xrightarrow{\text{PD}}$

If $P \in S_+^\gamma \Leftrightarrow P \succeq 0$

$P \in S_{++}^\gamma \Leftrightarrow P \succ 0$

If $P \in S_+^\gamma \Leftrightarrow x^T P x \geq 0 \quad \forall x \in \mathbb{R}^n$

If $P \in S_{++}^\gamma \Leftrightarrow x^T P x > 0 \quad \forall x \in \mathbb{R}^n$

If $P \in S_+^\gamma \Rightarrow$ Define P^{γ_2}

$$P^{\gamma_2} = U \Lambda^{\gamma_2} U^T$$

$$\hookrightarrow \text{diag}(\sqrt{\lambda_1}, \sqrt{\lambda_2}, \dots, \sqrt{\lambda_n})$$

$$P^{\gamma_2} P^{\gamma_2} = P$$

$$(\gamma_2, \dots, \gamma_2) (\gamma_1, \dots, \gamma_1)$$

$$(U \Lambda^{Y_2} U^\top) \underbrace{(U \Lambda^{Y_2} U^\top)}_{=I}$$

$$U \Lambda^{Y_2} \Lambda^{Y_2} U^\top = U \Lambda U^\top = P$$

P-Quadratic norm

Let $P \in S_+$

$$\|x\|_P = \underbrace{(x^\top P x)^{1/2}}_{P=I} = \underbrace{\|P^{1/2}x\|_2}_{(x^\top x)^{1/2}} = \|x\|_2$$

$$\Leftrightarrow \|x\|_2 = \|x\|_I$$

$$P > Q \Leftrightarrow P - Q > 0 \Rightarrow \text{Notation}$$

Norms on matrices $\Rightarrow X \in \mathbb{R}^{m \times n}$

$$\|X\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |X_{ij}|^2}$$

$$\|X\|_{\text{sum}} = \sum_{i=1}^m \sum_{j=1}^n |X_{ij}|$$

$$\|X\|_{\max} = \max \left\{ |X_{ij}|, i=1, \dots, m; j=1, \dots, n \right\}$$

All norms in finite-dimensional spaces (in particular, \mathbb{R}^n)

All norms in finite-dimensional spaces (in particular, \mathbb{R}^n) are "equivalent".

There always exists scalars α and β (which may depend on n) for given norms $\|\cdot\|_a$ and $\|\cdot\|_b$ s.t.

$$\alpha \|\cdot\|_a \leq \|\cdot\|_b \leq \beta \|\cdot\|_a$$

e.g., $\|\cdot\|_1 \leq n \|\cdot\|_\infty$

\downarrow

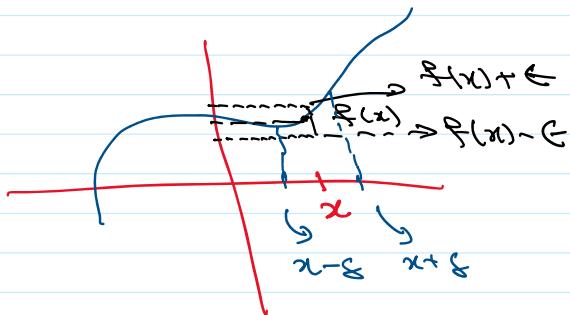
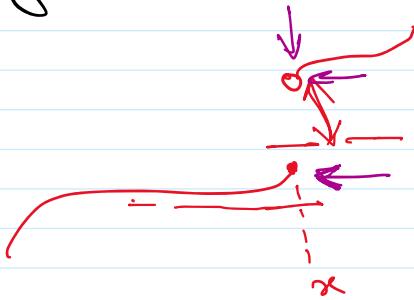
$\beta = n$

Calculus ConceptsContinuous functions

A function $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is continuous at $x \in \mathbb{R}^n$ if for any $\epsilon > 0$, $\exists \delta > 0$ such that

$$\|y-x\|_2 \leq \delta \Rightarrow \|f(y) - f(x)\|_2 \leq \epsilon$$

for all $y \in \text{dom } f$.



\Leftrightarrow Given any sequence x_1, x_2, \dots
Let $\lim_{n \rightarrow \infty} x_n = x$

* f is continuous at x if and only if

$$\lim_{n \rightarrow \infty} f(x_n) = f(\lim_{n \rightarrow \infty} x_n) = f(x)$$

f is a continuous function if it is continuous at all $x \in \text{dom } f$

Derivative of a function $f: \mathbb{R} \rightarrow \mathbb{R}$

Suppose $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $x \in \text{int domf}$

The function f is differentiable at $x \in \mathbb{R}^n$ if
 \exists a matrix $Df(x) \in \mathbb{R}^{mn}$ that satisfies:

$$\lim_{\substack{z \in \text{domf}, z \neq x \\ z \rightarrow x}} \frac{\|f(z) - f(x) - Df(x)(z-x)\|_2}{\|z-x\|_2} = 0$$

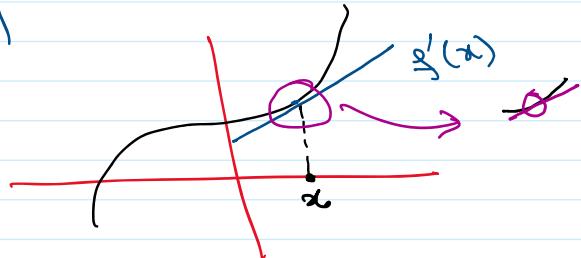
$Df(x)$ is called derivative of f at x

\hookrightarrow when $m > 1 \Rightarrow$ Jacobian of f .

Special case: $n = 1, m = 1$

$$\lim_{\substack{z \in \text{domf}, z \neq x \\ z \rightarrow x}} \frac{|f(z) - f(x) - f'(x)(z-x)|}{|z-x|} = 0$$

$$\lim_{\substack{z \in \text{domf}, z \neq x \\ z \rightarrow x}} \frac{|f(z) - f(x) - f'(x)(z-x)|}{|z-x|} = 0$$



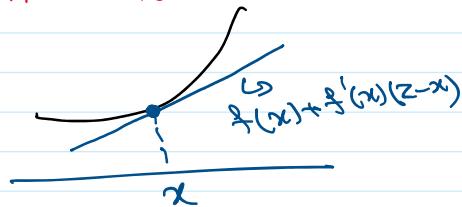
The function f is differentiable if domf is open and f is differentiable at every $x \in \text{domf}$.

Derivatives provide first-order (or linear) approximation of f at x

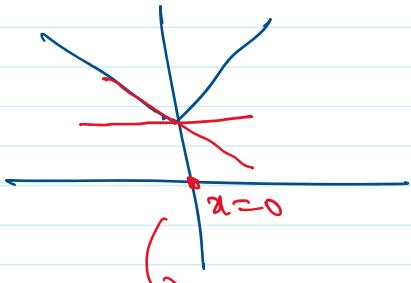
$$\hat{f}(z) = f(x) + f'(x)(z-x)$$

↳ linear function of z

Affine function



$Df(x)$ is unique.



first-order approximation of a function

Not differentiable at $x=0$

$$f(x) = 3x + 2$$

$$f'(x) = 3$$

$$\begin{aligned} \hat{f}(z) &= 3z + 2 + 3(z-2) \\ &= 3z + 2 + 3z - 6 \\ &= 3z + 2 \end{aligned}$$

$$\left[Df(x) \right]_{i,j} = \frac{\partial f_i(x)}{\partial x_j} ; \quad \begin{matrix} i=1, \dots, m \\ j=1, \dots, n \end{matrix}$$

$$f(x) : \mathbb{R}^2 \rightarrow \mathbb{R} ; \quad f(x) = x_1^2 + x_2^2$$

$$Df(x) : \mathbb{R}^{1 \times 2} ; \quad \left[Df(x) \right]_{1,1} = \frac{\partial f(x)}{\partial x_1} = 2x_1$$

$$\left[Df(x) \right]_{1,2} = \frac{\partial f(x)}{\partial x_2} = 2x_2$$

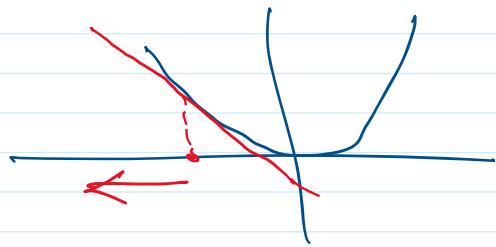
$$Df(x) = \begin{bmatrix} 2x_1 & 2x_2 \end{bmatrix}$$

when $m=1 \Rightarrow Df(x)$ is a row vector of length n

$Df(x)^\top = \nabla f(x) \Rightarrow$ Gradient of the function at x

$$[\nabla f(x)]_i = \frac{\partial f(x)}{\partial x_i}, i=1, \dots, n$$

$$f(x) = x^2$$



Scalar-valued functions

$$\hat{f}(z) = f(x) + \nabla f(x)^\top (z-x)$$

$$\hat{f}(z) \rightarrow f(z) \text{ as } z \rightarrow x$$

Example

$$\textcircled{1} \quad f(x) = \frac{1}{2} \underbrace{x^T P x}_{} + q^T x + r \quad ; \quad f: \mathbb{R}^n \rightarrow \mathbb{R}$$

$P \in \mathbb{S}^n$ and $q \in \mathbb{R}^n$ and $r \in \mathbb{R}$

$$g(x) = ax^2 + bx + c$$

$Df(x)$?

$$\nabla f(x) = Df(x)^T$$

$$Df(x) = 2 \cdot \frac{1}{2} x^T P + q^T$$

$$\nabla f(x) = Px + q$$

$$\textcircled{2} \quad f(x) = \log(\det x) ; \quad \det x > 0$$

$\hookrightarrow \text{dom } f = \mathbb{S}_{++} \Rightarrow$ Positive definite matrices

$Df(x)$ or $\nabla f(x)$?

Make use of the fact that $Df(x)$ provides a linear approximation of $f(z)$, when z is 'close' to x .

Consider $z \in \mathbb{S}_{++}$ such that z is close to x

$$\Rightarrow z = x + \Delta x, \text{ where } \Delta x \rightarrow 0$$

as $z \rightarrow x$

Note: We cannot claim that $\Delta x \in \mathbb{S}_{++}$

Goal: show that $f(z) = f(x) + \underbrace{Df(x)}_{\text{Derivative}}(z-x)$

as $z \rightarrow x \Rightarrow Df(x)$ would be our derivative

$$z = x + \Delta x$$

$$= (x^{\gamma_2} (I + x^{\gamma_2} \Delta x x^{\gamma_2}) x^{\gamma_2})$$

$$\log \det(z) = \log \det(x^{\gamma_2} (I + x^{\gamma_2} \Delta x x^{\gamma_2}) x^{\gamma_2})$$

$$\underbrace{f(z)}_{\log \det} = \log [(\det x^{\gamma_2}) (\det(I + x^{\gamma_2} \Delta x x^{\gamma_2})) (\det x^{\gamma_2})]$$

$$= \log \underbrace{[\det x^{\gamma_2} \cdot \det x^{\gamma_2} \cdot \det(I + x^{\gamma_2} \Delta x x^{\gamma_2})]}_{\det x}$$

$$= \underbrace{\log \det x}_{f(x)} + \underbrace{\log \det(I + x^{\gamma_2} \Delta x x^{\gamma_2})}_{\circlearrowleft (z-x)}$$

$I + x^{\gamma_2} \Delta x x^{\gamma_2}$ \Rightarrow This has eigenvalue decomposition

$\lambda_1, \lambda_2, \dots, \lambda_n$ are the eigenvalues of $x^{\gamma_2} \Delta x x^{\gamma_2}$

$$\text{eigenvalues of } (I + A) = 1 + \text{eigenvalues}(A)$$

$1 + \lambda_i, i=1, \dots, n$ are the eigenvalues of A

↗

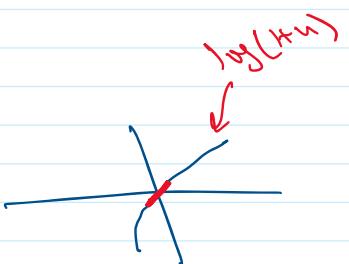
$1 + \lambda_i, i=1, \dots, r$ are the eigenvalues of
 $\det(I + X^{-1/2} \Delta X X^{-1/2})$
 $\hookrightarrow = \prod_{i=1}^r (1 + \lambda_i)$

$$f(z) = f(x) + \log \left(\prod_{i=1}^r (1 + \lambda_i) \right)$$

$$= f(x) + \sum_{i=1}^r \log(1 + \lambda_i).$$

Since $\Delta X \rightarrow 0 \Rightarrow X^{-1/2} \Delta X X^{-1/2} \rightarrow 0$

$$\Rightarrow \lambda_i \rightarrow 0 \quad \forall i$$



$$\Rightarrow \text{As } z \rightarrow x, \lambda_i \rightarrow 0$$

$$\log(1+u) \approx u \text{ for } u \text{ very small}$$

$$\Rightarrow \sum_{i=1}^r \log(1 + \lambda_i) \approx \sum_{i=1}^r \lambda_i \text{ when } z \rightarrow x$$

$$f(z) = f(x) + \sum_{i=1}^r \lambda_i \quad \text{as } z \rightarrow x$$

\downarrow
eigenvalues of $X^{-1/2} \Delta X X^{-1/2}$

$\sum \text{eigenvalues of } A = \text{tr}(A)$

$\Rightarrow = f(x) + \text{tr}(X^{-1/2} \Delta X X^{-1/2})$

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

$$= f(x) + \nabla f(x)^T \Delta x$$

$$= f(x) + \nabla f(x)^T \Delta x$$

$\hookrightarrow (z-x)$

$$f(z) = f(x) + \langle x^T, z-x \rangle$$

$$\nabla f(x) = x^T$$

Chain rule

It applies in higher dimensions also.

scalar: $h(u) = g(f(u))$

$$h'(u) = g'(f(u)) f'(u)$$

If $h(x) = g(f(x))$; $x \in \mathbb{R}^n$

$$Dh(x) = Dg(f(x)) \cdot D_x f(x)$$

Example: Say $g(x) = f(Ax+b)$

$$\begin{aligned} Dg(x) &= Df(Ax+b) - D(Ax+b) \\ &= Df(Ax+b) - A \end{aligned}$$

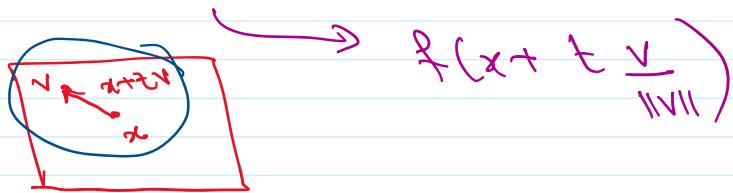
$$\nabla g(x) = A^T \nabla f(Ax+b)$$

Directional derivative of a function

Let $f: \mathbb{R}^n \rightarrow \mathbb{R}$, $x \in \mathbb{R}^n$

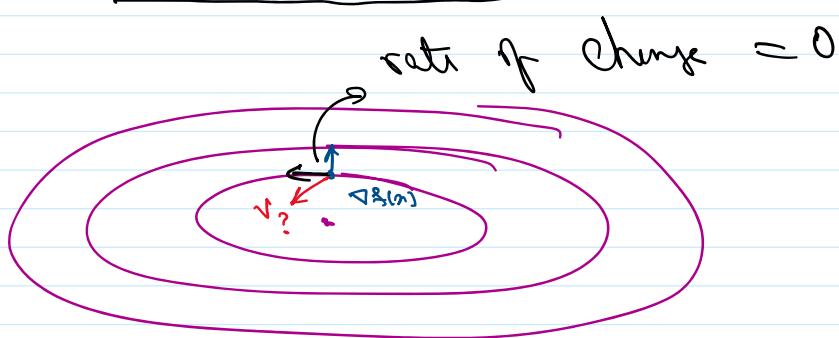
Directional derivative along a vector $v \in \mathbb{R}^n$ at $x \in \mathbb{R}^n$ is defined as follows:

$$\tilde{f}(t) = f(x + tv) : \mathbb{R} \rightarrow \mathbb{R}$$



$$\begin{aligned}\tilde{f}'(t) &= D\tilde{f}(x + tv) \cdot D(x + tv) \\ &= \nabla f(x + tv)^T \cdot v\end{aligned}$$

Directional derivative of $f(x)$ in the direction v is defined as $\tilde{f}'(0) = \nabla f(x)^T v$



Algorithms for unconstrained Optimization

$f: \mathbb{R}^n \rightarrow \mathbb{R}$; $\text{dom } f$

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

$$\min_{\underline{x} \in \text{dom}} f(x)$$

Assume the minimum value is attained by $f(x)$.

$$\text{Define } p^* = \min_x f(x)$$

An optimization algorithm is an iterative method that produces a sequence of points $x^{(0)}, x^{(1)}, \dots, x^{(k)}$ such that

$$f(x^{(k)}) \rightarrow p^* \text{ as } k \rightarrow \infty$$

In the case when $\arg \min_x f(x)$ is unique ($= x^*$)

then we also hope that $x^{(k)} \rightarrow x^*$ as $k \rightarrow \infty$

Suppose $f(\cdot)$ is continuous

$$f(x^{(k)}) \rightarrow f(x^*) = p^*$$

Question: When to terminate the algorithm?

maybe $f(x^{(k)}) - p^* \leq \epsilon$ for ϵ small \Rightarrow terminate

↳ we do not know $p^* \Rightarrow$ not practical solution.

Search Direction-based Iterative Optimization Algorithms

Pseudo Code

Initialize: $x^{(0)} \in \text{dom}$

$$k \leftarrow 0$$

while stopping criterium not satisfied

do

$$x^{(k+1)} \leftarrow x^{(k)} + t^{(k)} \Delta x^{(k)}$$

$$k \leftarrow k+1$$

$t^{(k)} \in \underbrace{\mathbb{R}_{++}}_{>0} \rightarrow$ Step size \rightarrow (learning rate in stochastic optimization / machine learning)

$\Delta x^{(k)}$ \rightarrow Search direction at time k

$x^{(k)}$ \Rightarrow Iterate

Main challenge: ① How to pick $\Delta x^{(k)}$?
② How to pick $t^{(k)}$?

Stopping Criterion

- Objective function stops changing significantly

$$|f(x^{(k)}) - f(x^{(k-1)})| \leq \epsilon \text{ for } \epsilon \text{ small}$$

E.g. $\epsilon = 10^{-8}$

- Iterates stop changing significantly

$$\|x^{(k)} - x^{(k-1)}\| \leq \epsilon \text{ for } \epsilon \text{ small}$$

- Function gradient evaluated at the iterates becomes very small

$$\|\nabla f(x^{(k)})\| \leq \epsilon \text{ for } \epsilon \text{ small}$$

Descent Optimization Methods

An optimization method is termed a 'descent method' if

$$f(x^{(k+1)}) < f(x^{(k)}) \quad \forall k, \text{ except when } x^{(k)} = x^*.$$

strict inequality

when we initialize at $x^{(0)}$

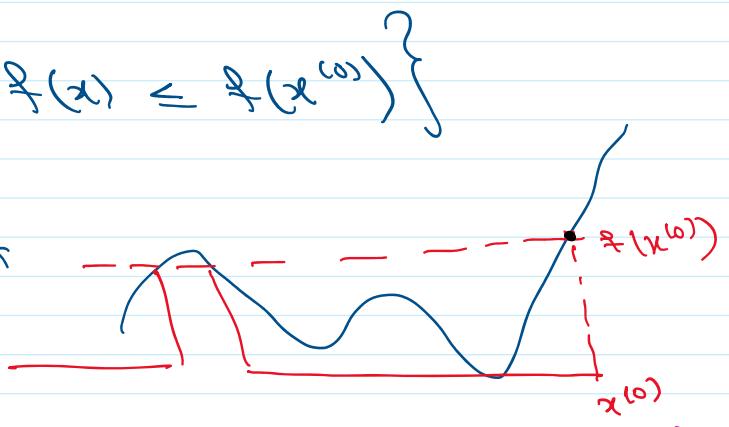
All iterates $x^{(k)}$ for a descent method stay within the set

$$\{ \dots, x_1, x_2, \dots, x_n \}$$

the set

$$S = \{x : f(x) \leq f(x^{(0)})\}$$

Sublevel set of $f(x)$ at
 $f(x) = f(x^{(0)})$



$$x^{(w+1)} \leftarrow x^{(w)} + t^{(w)} \Delta x^{(w)}$$

In descent methods, $\Delta x^{(w)}$ is called descent direction.

What direction is a descent direction?

We will focus on all continuously differentiable functions and answer this question.

$$C' = \{f(x) : f(x) \text{ has continuous derivatives at all } x \in \text{dom} f\}$$

① Make use of Taylor's theorem

Since $f \in C'$, we have that

$$f(z) = f(x) + \nabla f(x)^T (z-x) + h(x) \|z-x\|_2$$

where $h(x) \rightarrow 0$ as $z \rightarrow x$ faster than $\|z-x\|_2$

$$f(z) = f(x) + \nabla f(x)^T (z-x) + o(\|z-x\|_2)$$

$$\text{P.I. } 0.1n - \overset{-(2+\epsilon)}{\sim}$$

$$\text{e.g. } g(n) = n^{-(2+\epsilon)}$$

$$g(n) = o(n^{-2})$$

$$\frac{h(x) \|z-x\|_2}{\|z-x\|_2} \xrightarrow{\quad} 0$$

Take

$$x = x^{(u)}$$

$$z = x^{(u+1)} = x^{(u)} + t^{(u)} \Delta x^{(u)}$$

$$f(z^{(u+1)}) = f(x^{(u)}) + \nabla f(x^{(u)})^T (t^{(u)} \Delta x^{(u)}) + o(\|t^{(u)} \Delta x^{(u)}\|_2)$$

$$= f(x^{(u)}) + t^{(u)} \nabla f(x^{(u)})^T \Delta x^{(u)} + o(t^{(u)} \|\Delta x^{(u)}\|_2)$$

Let $t^{(u)} \rightarrow 0$ (very small)

$$f(z^{(u+1)}) = f(x^{(u)}) + t^{(u)} \nabla f(x^{(u)})^T \Delta x^{(u)}$$

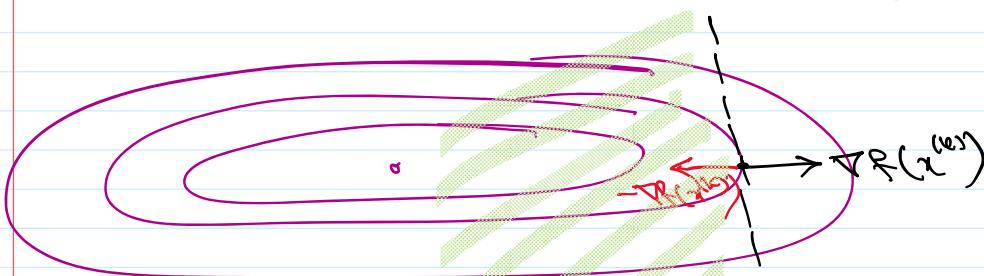
$$f(z^{(u+1)}) < f(x^{(u)})$$

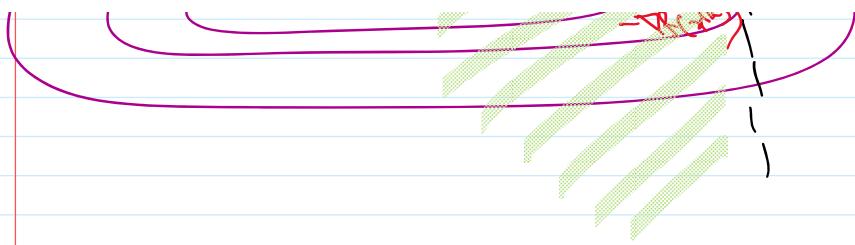
$$\Leftrightarrow \nabla f(x^{(u)})^T \Delta x^{(u)} < 0$$

$$\Leftrightarrow [-\nabla f(x^{(u)})]^T \Delta x^{(u)} > 0$$

$\Delta x^{(u)}$ is a descent if and only if

It makes an acute angle with $-\nabla f(x^{(u)})$





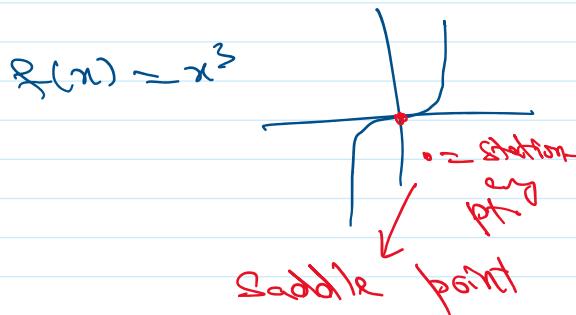
In particular, $\Delta x^{(k)} = -\nabla f(x^{(k)})$ is a descent direction

In a descent method, the optimization method might not be able to reduce function value further when

$$\nabla f(x^{(k)}) = 0$$

Strict use of a descent method means that any x for which $\nabla f(x^{(k)}) = 0$ is a fixed point of the method.

Any $x \in \text{dom } f$ for which $\nabla f(x) = 0$ is called a stationary point of f .



General form of descent direction

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

$$\boxed{\Delta x^{(k)} = -B^{(k)} \nabla f(x^{(k)})}$$



Where $B^{(k)}$ is a positive definite matrix; $\in S_{++}^n$

$$\begin{aligned}
 \text{Descent methods} \Rightarrow \boldsymbol{x}^{(k+1)} &= \boldsymbol{x}^{(k)} - t^{(k)} \boldsymbol{B}^{(k)} \nabla f(\boldsymbol{x}^{(k)}) \\
 &\left[-\nabla f(\boldsymbol{x}^{(k)}) \right]^\top \Delta \boldsymbol{x}^{(k)} \\
 &= -\nabla f(\boldsymbol{x}^{(k)})^\top (-\boldsymbol{B}^{(k)} \nabla f(\boldsymbol{x}^{(k)})) \\
 &= \nabla f(\boldsymbol{x}^{(k)})^\top \boldsymbol{B}^{(k)} \nabla f(\boldsymbol{x}^{(k)}) \\
 &> 0 \text{ b/c } \boldsymbol{B}^{(k)} \text{ is PD.}
 \end{aligned}$$

Based on the choice of $\boldsymbol{B}^{(k)}$, we have different names for descent methods.

$$\begin{aligned}
 \textcircled{1} \text{ Gradient descent : } \boldsymbol{B}^{(k)} &= \mathbf{I} \\
 \Leftrightarrow \Delta \boldsymbol{x}^{(k)} &= -\nabla f(\boldsymbol{x}^{(k)})
 \end{aligned}$$

$$\boldsymbol{x}^{(k+1)} = \boldsymbol{x}^{(k)} - t^{(k)} \nabla f(\boldsymbol{x}^{(k)})$$

\textcircled{2} Newton's method

$$\begin{aligned}
 \boldsymbol{B}^{(k)} &= \text{Hessian matrix at } \boldsymbol{x}^{(k)} \\
 &= \nabla^2 f(\boldsymbol{x}^{(k)})
 \end{aligned}$$

Hessian matrix \Rightarrow matrix of second-order partial derivatives.

Another type is called Quasi-Newton method, in which $\boldsymbol{B}^{(k)}$ is built from $\nabla f(\boldsymbol{x}^{(k)})$, but is meant to be \mathbf{I} .

in which $B^{(k)}$ is built from $\nabla^2 f(x^{(k)})$, but is meant to approximate $\nabla^2 f(x^{(k)})$

③ Steepest descent

$B^{(k)}$ is chosen based on the geometry of the function.

Issue: The previous analysis guarantees a descent direction, but only when $t^{(k)}$ is very small.

Can we use descent methods with a larger step size?

↳ we can, but we need to assume additional regularity on the function.

$$C_2(R^n) = \{ f(x) : f \text{ has continuous derivatives that are } L\text{-Lipschitz continuous} \}$$

$C_2 \subseteq C$ $\Rightarrow L\text{-Smooth functions}$

The gradients $\nabla f(x)$ of $f(x)$ are called L -Lipschitz continuous if and only if

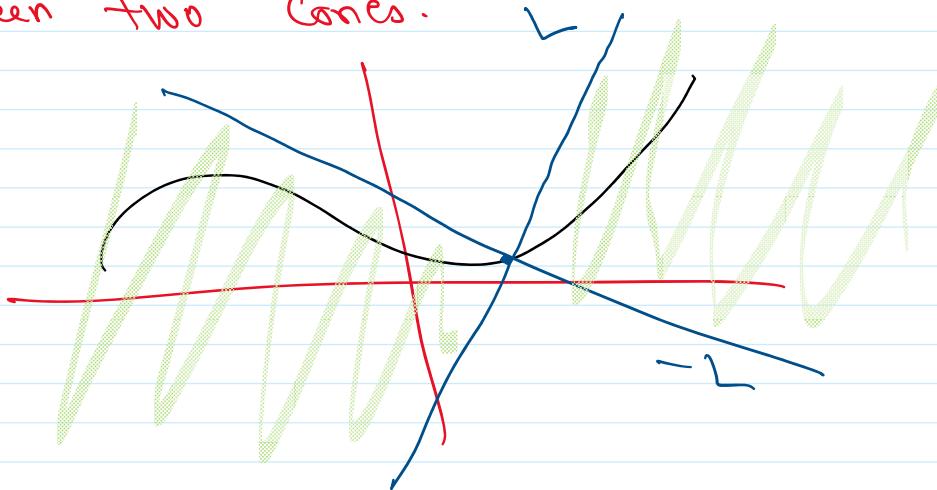
$$\| \nabla f(x) - \nabla f(y) \|_2 \leq L \|x-y\|_2, \forall x, y$$

defn

\downarrow

Lipschitz Constant

Lipschitz Continuity is a stronger form of continuity, which says that the function value always lies between two cones.



e.g. ① $f(x) = x^2$

$$f'(x) = 2x$$

$$|f'(x) - f'(y)| = |2x - 2y| \leq 2|x-y|$$

$$L = 2$$

② $f(x) = x_1^2 + x_2^2$

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$\nabla f(x) = 2x = \begin{bmatrix} 2x_1 \\ 2x_2 \end{bmatrix}$$

$$\|\nabla f(x) - \nabla f(y)\|_2 = \|2x - 2y\|_2 \leq 2\|x - y\|_2$$

$$L = 2$$

$f \in C_L^1(\mathbb{R}^n) \Rightarrow$ Continuously differentiable functions with Lipschitz gradients

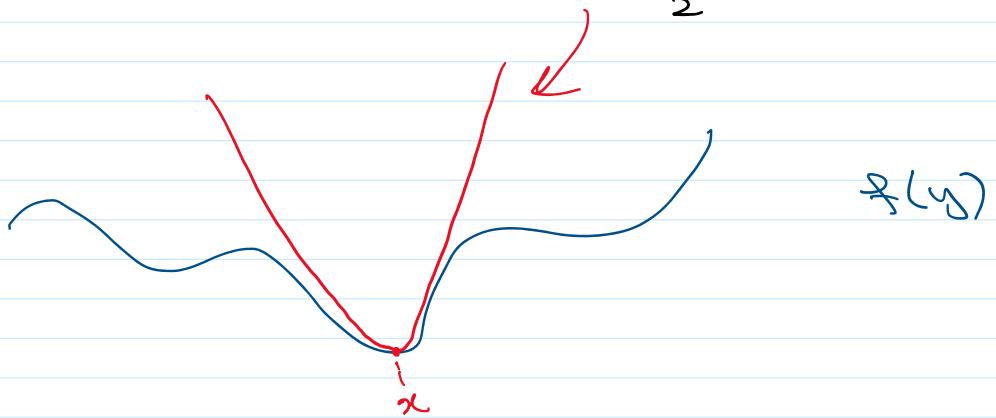
Lemma: Quadratic upper bound on $f \in C_L^1(\mathbb{R}^n)$

Let $f \in C_L^1(\mathbb{R}^n)$ and dom f is a convex set.

Then; $\forall x, y \in \text{dom } f$

→ quadratic function in y

$$f(y) \leq f(x) + \nabla f(x)^T (y-x) + \frac{\|y-x\|_2^2}$$



Proof:

Define $g(t) = f(x + t(y-x))$ for $t \in [0,1]$

since dom f is convex $(x + t(y-x))$ is in the dom f.

$$\Rightarrow x + ty - tx = (1-t)x + t y$$

Note: $g(1) = f(y)$

$g(0) = f(x)$

$$g'(t) = \nabla f(x + t(y-x))^T (y-x)$$

$$\nabla g'(0) = \nabla f(x)^T (y-x)$$

$$g'(0) = \nabla f(x)^T (y-x)$$

$$\begin{aligned} g'(t) - g'(0) &= \nabla f(x + t(y-x))^T (y-x) - \nabla f(x)^T (y-x) \\ &= [\nabla f(x + t(y-x)) - \nabla f(x)] (y-x) \end{aligned}$$

Use Cauchy-Schwarz Inequality ($a^T b \leq \|a\|_2 \|b\|_2$)

$$\leq \|\nabla f(x + t(y-x)) - \nabla f(x)\|_2 \|y-x\|_2$$

Lipschitz continuous gradients

$$\leq L \|x + t(y-x) - x\|_2 \|y-x\|_2$$

$$g'(t) - g'(0) \leq t L \|y-x\|_2^2$$

$$\int_0^1 g'(t) dt = g(1) - g(0)$$

$$\Rightarrow g(1) = g(0) + \int_0^1 g'(t) dt$$

$$\leq g(0) + \int_0^1 (t L \|y-x\|_2^2 + g'(0)) dt$$

$$f(y) \leq f(x) + L \|y-x\|_2^2 \left. \frac{t^2}{2} \right|_0^1 + g'(0)$$

$$\leq f(x) + \nabla f(x)^T (y-x) + \frac{L}{2} \|y-x\|_2^2$$



Descent Lemma (when $\lambda x = \nabla f(x)$)

CC

Descent Lemma (when $\Delta x = \nabla f(x)$)

Let $f \in C_1(\mathbb{R}^n)$ with $\text{dom } f$ being convex. Let us consider the iterative method:

$$x^{(k+1)} = x^{(k)} + t^{(k)} \Delta x^{(k)}$$

and focus on the case $\Delta x^{(k)} = -\nabla f(x^{(k)})$. Then, as long as $0 < t^{(k)} < \frac{2}{L}$, $\Delta x^{(k)}$ is a descent direction, i.e. $f(x^{(k+1)}) < f(x^{(k)})$ as long as $x^{(k)}$ is not x^* .

Remarks: A range of step sizes guarantee descent. The larger L is, the smaller is the range.

Proof: Let $y = x^{(k)} + t^{(k)} \Delta x^{(k)}$

and $\alpha = x^{(k)}$ in the Quadratic Upper bound.

$$\begin{aligned} f(x^{(k)} + t^{(k)} \Delta x^{(k)}) &\leq f(x^{(k)}) + \nabla f(x^{(k)})^\top (x^{(k)} + t^{(k)} \Delta x^{(k)} - \alpha) \\ &\quad + \frac{1}{2} \|x^{(k)} + t^{(k)} \Delta x^{(k)} - x^{(k)}\|_2^2 \end{aligned}$$

$$f(x^{(k+1)}) \leq f(x^{(k)}) + t^{(k)} \nabla f(x^{(k)})^\top \Delta x^{(k)} + t^{(k)} \frac{L}{2} \|\Delta x^{(k)}\|_2^2$$

Now put $\Delta x^{(k)} = -\nabla f(x^{(k)})$

$$f(x^{(k+1)}) \leq f(x^{(k)}) - t^{(k)} \|\nabla f(x^{(k)})\|_2^2 + t^{(k)} \frac{L}{2} \|\nabla f(x^{(k)})\|_2^2$$

$$\text{need} < 0$$

$$\leq f(x^{(k)}) - \left(t^{(k)} - \frac{t^{(k)2}}{L}\right) \|\nabla f(x)\|_2^2$$

Clearly, $f(x^{(k+1)}) < f(x^{(k)})$

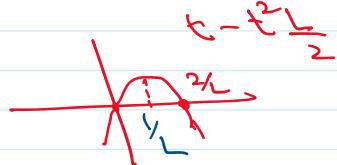
if and only if

$$t^{(k)} - \frac{t^{(k)2}}{L} > 0$$

$\Leftrightarrow \frac{t^{(k)2}}{L} < t^{(k)}$

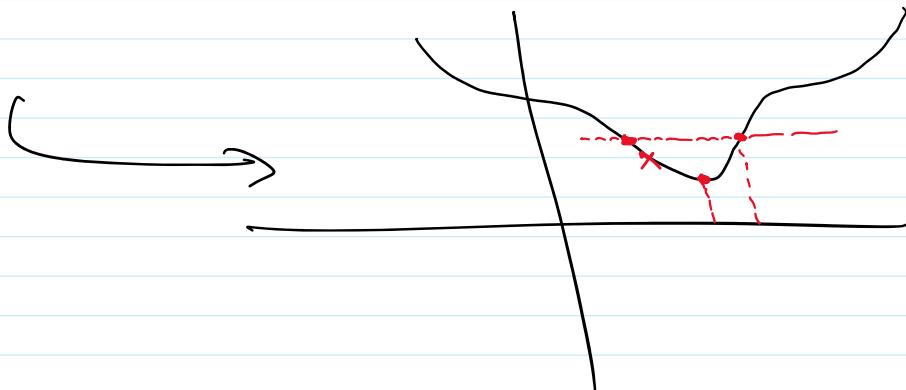
$\Leftrightarrow t^{(k)} < \frac{2}{L}$

What is the $t^{(k)}$ when ' L ' is known?



The best $t^{(k)}$, when it comes to most reduction for gradient descent is $t^{(k)} = \frac{2}{L}$.

Ex: $f(x)$



General Descent Method

Initialize: $x^{(0)} \in \text{dom}$
 $k \leftarrow 0$

Repeat

1. Determine a descent direction $\Delta x^{(k)}$
(i.e., $-\nabla f(x^{(k)})^\top \Delta x^{(k)} > 0$)

2. Line Search: choose a step size $t^{(k)} > 0$

3. Update the iterate: $x^{(k+1)} \leftarrow x^{(k)} + t^{(k)} \Delta x^{(k)}$

Until Stopping criterion is satisfied.

Special Cases:

① Gradient descent: $\Delta x^{(k)} = -\nabla f(x^{(k)})$

② Newton's method: $\Delta x^{(k)} = -\nabla^2 f(x^{(k)})^{-1} \nabla f(x^{(k)})$

Line Search: How to pick the step size in a descent method?

① Choose a fixed step size $t^{(k)} = \eta$, $\forall k \geq 0$.

e.g.) If $f \in C_1(\mathbb{R}^n)$ and L is known or can be computed efficiently then pick

$\eta = \frac{1}{L}$ for gradient descent.

In other cases, trial and error helps pick a step size.

Issues with this approach:

- Sometimes it is too costly to compute L .
- Sometimes functions are not C^1 .
- Even when one has descent in some iterations, does not mean the step size would work in all points in domf.

Unless L is known, the larger the step size, the better 'perhaps'.

② Variable Step size $t^{(k)}$

↳ The most common approach in the literature.

(a) Decaying step size policy

$$t^{(k)} \rightarrow 0 \text{ as } k \rightarrow \infty$$

Typical policy

$$(i) t^{(k)} \rightarrow 0 \text{ as } k \rightarrow \infty$$

$$(ii) \sum_{k=0}^{\infty} t^{(k)} = \infty$$

$$(iii) \sum_{k=0}^{\infty} [t^{(k)}]^2 < \infty$$

e.g., $t^{(k)} = \frac{\text{const}}{k}$

↳ Often used in machine learning / stochastic

↳ Often used in machine learning / Stochastic optimization.

(b) Search for a nice step size that reduces the objective function using a subroutine in each iteration k .

↳ often used in practical deterministic optimization problems.

↳ we will study this in detail under 'Exact line search' and 'Inexact line search'.

Exact line search

$$x^{(k+1)} = x^{(k)} + t^{(k)} \Delta x^{(k)}$$

When we have fixed $\Delta x^{(k)}$ and are looking for $t^{(k)}$, we are effectively looking at a one-dimensional function:

$$\tilde{f}(t) = f(x^{(k)} + t \Delta x^{(k)})$$

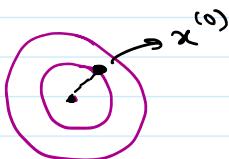
$$\text{Pick } t^{(k)} = \arg \min_{t \geq 0} \tilde{f}(t)$$

Example: say $f(x) = x_1^2 + x_2^2$

$$\nabla f(x) = 2x$$

$$\text{Let } x^{(0)} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \mathbf{1}$$

$$\Rightarrow x^{(1)} = x^{(0)} - 2t^{(1)} x^{(0)}$$



$$\begin{aligned}
 &= (1-2t^{(0)}) \frac{1}{2} = \left(1 - \frac{1}{2} \times 2\right) \frac{1}{2} = 0. \frac{1}{2} \\
 \hat{f}(t) &= f(x^{(0)}) = (1-2t)^2 \cdot 1 + (1-2t)^2 \cdot 1 \\
 t^{(0)} &\approx \gamma_2 = \arg \min_{t \geq 0} \hat{f}(t) = \arg \min_{t \geq 0} 2(1-2t)^2
 \end{aligned}$$

$$\hat{f}'(t) = 4(1-2t) \cdot -2 = -8(1-2t)$$

$$\hat{f}'(t) = 0 \Rightarrow -8(1-2t) = 0 \\ t = \gamma_2$$

Exact line search, in which one solves a one-dim optimization problem in each iteration, works in cases where:

- ① The solution has an analytical form.
- ② It might be computationally feasible to numerically solve the 1-D problem.

But if the cost of exact line search is too much, we resort to inexact line search.

↳ Backtracking (Armijo-Goldstein line search)

$f \in C_L(\mathbb{R}^n) \Rightarrow f$ is L -smooth

Function is Smooth (or non-smooth)

Quadratic Function

$$\underline{f(x)} = \frac{1}{2} x^T Q x + b^T x + r$$

Generally, we focus on $f(x)$ being convex

$\Rightarrow Q$ is Positive Semi definite or Positive definite

$$f(x) = a x^2$$

\hookleftarrow Convex when $a > 0$

Special case: $f(x) = \frac{1}{2} x^T Q x$; Q is P.D.
 $Q \in \mathbb{S}_{++}$

$f(x)$ in this case is L -smooth

Linear algebra review

(Appendix A.1.5)

Operator norm of a matrix

\hookrightarrow Singular Value Decomposition of a matrix

Every matrix $A \in \mathbb{R}^{m \times n}$ has a Singular Value decomposition

$$A = U \sum V^T$$

↓ ↓ ↓

$m \times m$ $n \times n$ matrix with
orthonormal columns
(right singular vectors.)

↓
 $m \times m$
 with orthonormal
 columns

↪ left singular
 vectors of A

Diagonal matrix
 $m \times n$

"di"
 Singular
 Vectors.

singular values
 ≥ 0

$$A v_i = \sigma_i u_i$$

$$2 \times 3 \Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 \\ 0 & \sigma_2 & 0 \end{bmatrix}$$

$$\Sigma^T \Sigma = \begin{bmatrix} \sigma_1^2 & 0 & 0 \\ 0 & \sigma_2^2 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\begin{aligned}
 A A^T &= (U \Sigma V^T)(V \Sigma^T U^T) \\
 &= U \underbrace{\Sigma \Sigma^T}_{(m \times n)(n \times m)} V^T
 \end{aligned}$$

symmetric

entries are
 σ_i^2
 $U \Lambda U^T$
 $(m \times m)$

Eigenvalue decomposition
 of $A A^T$

① Singular values of A are the square root eigenvalues of $A A^T$ / $A^T A$

② Left singular vectors of A are the eigenvectors of $A A^T$

$$A^T A = \underbrace{U \Sigma^T \Sigma}_{{n \times n}} \underbrace{V^T V}_{{n \times n}} \underbrace{\Sigma^T \Sigma}_{{n \times n}}$$

entries are σ_i^2

③ Right singular vectors of A are the eigenvectors of $A^T A$

The operator norm of a matrix A with respect to the Euclidean norm is defined as:

$$\|A\|_2 = \sup \left\{ \|Ax\|_2 : \|x\|_2 \leq 1 \right\}$$

$$\begin{aligned} \|A\|_2 &= \sigma_{\max}(A) \Rightarrow \text{maximum Singular Value of } A \\ &= \sqrt{\lambda_{\max}(A^T A)} = \sqrt{\lambda_{\max}(AA^T)} \\ &\quad \text{maximum eigen value of } A^T A \end{aligned}$$

Submultiplicative property

$$\|Az\|_2 \leq \|A\|_2 \|z\|_2$$

General Definition of Operator Norm

$$\|A\|_{a,b} = \sup \left\{ \|Ax\|_a : \|x\|_b \leq 1 \right\}$$

Special Cases: $a=b$

$$\textcircled{1} \quad a=2$$

$$\textcircled{2} \quad a=1$$

$$\textcircled{3} \quad a=\infty$$

$$\|A\|_1 = \max_{j=1, \dots, n} \sum_{i=1}^m |A_{ij}|$$

$$\|A\|_\infty = \max_{i=1, \dots, m} \sum_{j=1}^n |A_{ij}|$$

Problem: Let $f(x) = \frac{1}{2} x^T Q x$ with $Q \in \mathbb{S}^n$.

Might not be Convex

Prove that $f(x)$ is L-smooth and derive the

Prove that $f(x)$ is L -smooth and derive the value of L . $\underbrace{f \in C^1_L}$

Solution:

$$\nabla f(x) = Qx$$

Let x and $y \in \mathbb{R}^n$

$$\|\nabla f(x) - \nabla f(y)\|_2 = \|Qx - Qy\|_2$$

$$= \|Q(x-y)\|_2$$

$$\leq \|Q\|_2 \|x-y\|_2$$

$$L = \|Q\|_2 = \sqrt{\lambda_{\max}(Q^\top Q)}$$

$$= \sqrt{\lambda_{\max}(Q^2)} = \sqrt{\lambda_{\max}(Q^2)}$$

$$= \lambda_{\max}(Q)$$

A with EVD $A = V \Lambda V^\top$

$$A^n = V \Lambda^n V^\top$$



Does Gradient Descent Converge?

If it does, at what rate?

Assume $f \in C^1_L(\mathbb{R}^n)$ and $b^* = \operatorname{argmin} f(x)$ exists.





$$f(y) \leq f(x) + \nabla f(x)^T (y-x) + \frac{L}{2} \|y-x\|_2^2$$

Put $y = x^{(k+1)} = x^{(k)} - t^{(k)} \nabla f(x^{(k)})$

$$x \approx x^{(k)}$$

$$t^{(k)} = \frac{1}{L}$$

$$f(x^{(k+1)}) \leq f(x^{(k)}) - \frac{1}{2} \|\nabla f(x^{(k)})\|_2^2 + \frac{1}{2L} \|\nabla f(x^{(k)})\|_2^2$$

$$f(x^{(k+1)}) \leq f(x^{(k)}) - \frac{1}{2L} \|\nabla f(x^{(k)})\|_2^2$$

↳ Descent lemma with $t^{(k)} = \frac{1}{L}$ and gradient descent

 $f(x^{(k+1)}) - f(x^{(k)}) \leq -\frac{1}{2L} \|\nabla f(x^{(k)})\|_2^2$

Add the above expression for $k=0, 1, \dots, n$

$$\sum_{k=0}^n [f(x^{(k+1)}) - f(x^{(k)})] \leq -\frac{1}{2L} \sum_{k=0}^{\infty} \|\nabla f(x^{(k)})\|_2^2$$

Telescoping sum $\lim_{k \rightarrow \infty} f(x^{(k+1)}) = f(x^{(\infty)})$

$$\sum_{n=0}^{\infty} (a_{n+1} - a_n) \Rightarrow \text{telescoping sum}$$

$$= \lim_{N \rightarrow \infty} \sum_{n=0}^N (a_{n+1} - a_n)$$

$$(a_{\underline{N+1}} - a_N) + (\cancel{a_N} - \cancel{a_{N-1}}) + (\cancel{a_{N-1}} - \cancel{a_{N-2}}) \\ + \dots + (\cancel{a_1} - \cancel{a_0})$$

$$= \lim_{N \rightarrow \infty} (a_{N+1} - a_0) = a_\infty - a_0$$

if $\lim_{N \rightarrow \infty} a_N = a_\infty$

$$\Rightarrow f(x^{(\infty)}) - f(x^{(0)}) \leq -\frac{1}{2L} \sum_{k=0}^{\infty} \|\nabla f(x^{(k)})\|_2^2$$

$$\boxed{\sum_{k=0}^{\infty} \|\nabla f(x^{(k)})\|_2^2} \leq 2L \cdot \underbrace{\left[f(x^{(0)}) - f(x^{(\infty)}) \right]}_{\geq p^*}$$

$$\leq 2L(f(x^{(0)}) - p^*)$$

$$\text{B/C } \sum_{k=0}^{\infty} \|\nabla f(x^{(k)})\|_2^2 < \infty$$

$$\Rightarrow \lim_{k \rightarrow \infty} \|\nabla f(x^{(k)})\|_2^2 = 0 \Leftrightarrow \lim_{k \rightarrow \infty} \nabla f(x^{(k)}) = 0$$

↳ we have converged to a local minimum or a saddle point.

↳ Basically a first-order stationary point.

When $f \in C^1(\mathbb{R}^n)$ and $t^{(k)} = \frac{1}{k}$

$$\|\nabla f(x^{(k)})\|_2^2 \rightarrow 0 \text{ as } k \rightarrow \infty$$

$$\nabla f(x^{(k)}) \rightarrow 0$$

What about the case of variable step sizes?

① Decaying step size policy

② Step size is bounded below;

Let $\epsilon > 0$ be a fixed constant

$$\frac{\epsilon}{2} \leq t^{(k)} \leq \frac{2-\epsilon}{2}$$

↳ Same proof works.

What about the rate of convergence?

$f \in C^1(\mathbb{R}^n)$; $t^{(k)} = \frac{1}{k}$

↳ From the previous lecture:

$$\star \quad \|\nabla f(x^{(k)})\|_2^2 \leq 2L [f(x^{(k)}) - f(x^{(k+1)})]$$

Sum \star from $k=1$ to $k=K$

$$\sum_{k=1}^K \|\nabla f(x^{(k)})\|_2^2 \leq 2L \sum_{k=1}^K [f(x^{(k)}) - f(x^{(k+1)})]$$

$$\leq 2L \sum_{k=0}^K \left\{ f(x^{(k)}) - f(x^{(k+1)}) \right\} > 0$$

↳ Telescoping sum

$$\leq 2L (f(x^{(0)}) - \underbrace{f(x^{(K+1)})}_{\geq f^*})$$

 $\sum_{k=1}^K \|\nabla f(x^{(k)})\|_2^2 \leq 2L (f(x^{(0)}) - f^*)$

$$\sum_{k=1}^K \|\nabla f(x^{(k)})\|_2^2 \geq K \min_{k \in \{1, 2, \dots, K\}} \|\nabla f(x^{(k)})\|_2^2$$

$$K \min_{k \in \{1, \dots, K\}} \|\nabla f(x^{(k)})\|_2^2 \leq 2L (f(x^{(0)}) - f^*)$$

$\gamma > 0$

$$\min_{k \in \{1, \dots, K\}} \|\nabla f(x^{(k)})\|_2^2 \leq \frac{\gamma}{K}$$

Within K iterations, we will have at least one $x^{(k)}$ such that $\|\nabla f(x^{(k)})\|_2^2 \leq \boxed{\frac{\gamma}{K}} = O(\frac{1}{K})$

Suppose we want $\|\nabla f(x^{(k)})\|_2^2 \leq \epsilon$ for ϵ very small

$$\Rightarrow \frac{\gamma}{K} \leq \epsilon \Rightarrow K \geq \frac{\gamma}{\epsilon}$$

$$\Rightarrow K = \Omega(\epsilon^{-1})$$

Say $\epsilon = 10^{-8}$

$\Rightarrow K = O(10^8)$ iterations

Similar results hold for step size choices for general descent methods, where the step size depends on the descent direction and is strictly lower bounded by $\epsilon > 0$.

Another Interpretation of Gradient Descent for $f \in C_1(\mathbb{R}^n)$

$$f(y) \leq f(x) + \nabla f(x)^T (y-x) + \frac{L}{2} \|y-x\|_2^2$$

Let us derive an iteration in which x is current iterate and $y=x^+$ is the next iterate

$$f(x^+) \leq f(x) + \nabla f(x)^T (x^+-x) + \frac{L}{2} \|x^+-x\|_2^2$$

We need x^+ such that $f(x^+)$ is as small as possible.

Mirror Descent proximity term
(when this is replaced by another $\|\cdot\|_1$)

We approach this by minimizing the upper bound

$$f(x^+) = f(x) + \nabla f(x)^T (x^+-x) + \underbrace{\frac{L}{2} \|x^+-x\|_2^2}_{\text{w.r.t. } x^+}$$

$$\begin{aligned} \nabla_{x^+} f(x^+) &= 0 + \nabla f(x)^T + \\ &\quad + \underbrace{\frac{L}{2} (2x^+ - 2x + 0)}_{\text{w.r.t. } x^+} \end{aligned}$$

$$\begin{aligned} (x^+-x)^T (x^+-x) &= x^{+T} x^+ - x^{+T} x \\ &\quad - x^{+T} x + x^{+T} x \\ &= x^{+T} x^+ - 2x^{+T} x \\ &\quad + x^{+T} x \end{aligned}$$

$$= \nabla f(x) + L(x^* - x) = 0$$

$$L(x^* - x) = -\nabla f(x)$$

$$\boxed{x^* = x - \frac{1}{L} \nabla f(x)}$$

Stepsize Selection when $f \in C_1(\mathbb{R}^n)$ but L is not known or computing it is too expensive.

↓
Inexact line search.

Backtracking | Armijo rule | Armijo-Goldstein step

Another approach based on Wolfe conditions, but they are harder to compute and we won't study them.

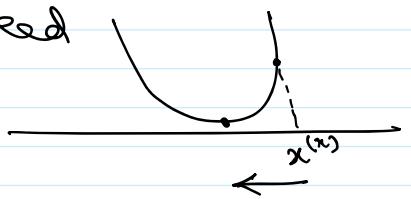
$$f^*(t) = f(x^{(u)} + t \Delta x^{(u)}) ; \quad t \geq 0 \\ t=0 \Rightarrow f(x^{(u)})$$

Inexact search line search requires finding a value of $t^{(u)}$ such that

$$f^*(t^{(u)}) = f(x^{(u)} + t^{(u)} \Delta x^{(u)})$$

sufficiently smaller than $f(x^{(u)})$

but there has to be a guaranteed decrease.



Algorithm (Backtracking)

Input: Current iterate x
Search direction Δx

Parameters $\alpha \in (0, 0.5) \rightarrow$ Sufficient decrease
 $\beta \in (0, 1)$ parameter

Initialize: $t \leftarrow 1$

while $f(x + t \Delta x) > f(x) + \alpha t \nabla f(x)^T \Delta x$

$$t \leftarrow \beta t$$

Sufficient decrease condition

The algorithm ends when $f(x + t \Delta x) \leq f(x) + \alpha t \nabla f(x)^T \Delta x$

$$\leftarrow 0$$

\hookrightarrow It depends on α .

Back tracking

$$t = 1$$

$$t = \beta$$

$$t = \beta^2$$

$\beta \Rightarrow$ The gridding of $(0, 1]$

\hookrightarrow larger β can slow down the line search
smaller β can end up giving you a very

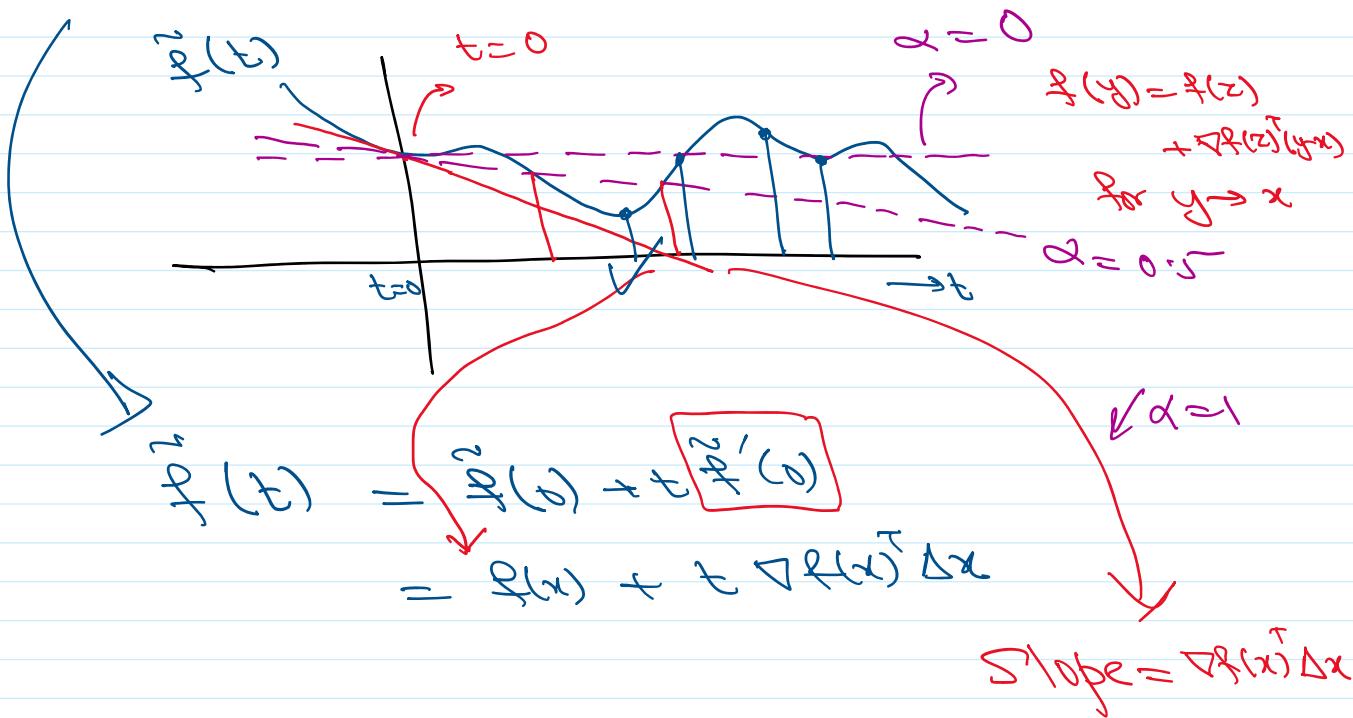
Small step size.

Geometric View of Backtracking

for t small enough

$$\hat{f}^2(t) = f(x + t\Delta x) \approx f(x) + t\nabla f(x)^T \Delta x$$

$$\hat{f}'(t) = \nabla f(x)^T \Delta x$$



what is the approximation when slope
is $\alpha \nabla f(x)^T \Delta x$

$$\hat{f}(t) = f(x) + \alpha t \nabla f(x)^T \Delta x$$

Switching to Newton's Method

Second Derivative of a function $f: \mathbb{R}^n \rightarrow \mathbb{R}$.

The second derivative of f , called the Hessian of f , at $x \in \text{int dom } f$, is denoted by $\nabla^2 f(x)$.

f , at $x \in \text{int dom } f$, is denoted by $\nabla^2 f(x)$
and is defined as:

$$\left[\nabla^2 f(x) \right]_{ij} = \frac{\partial^2 f(x)}{\partial x_i \partial x_j}, \quad i=1, \dots, n \\ j=1, \dots, n$$

$n \times n$
matrix

provided f is twice differentiable at x .

Gradient at $x \Rightarrow f(z) \approx f(x) + \nabla f(x)^T (z-x)$
 $as z \rightarrow x$

Hessian, by definition, is a quadratic approximation
of f at x

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

$\overbrace{\qquad\qquad\qquad}^{\hat{f}(z)}$

$\lim_{\substack{z \in \text{dom} \\ z \neq x \\ z \rightarrow x}} \frac{|f(z) - \hat{f}(z)|}{\|z-x\|^2} = 0$

Note: $D \nabla f(x) = \nabla^2 f(x)$

$\nabla f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^n$

\curvearrowleft Hessian is the derivative of the gradient

Next item is the Gradient of the

gradient

Chain rule for Hessians

① Suppose $f: \mathbb{R}^n \rightarrow \mathbb{R}$, $g: \mathbb{R} \rightarrow \mathbb{R}$

$$h(x) = g(f(x))$$

$$\nabla^2 h(x) = g'(f(x)) \nabla^2 f(x) + g''(f(x)) \nabla f(x) \nabla f(x)^T$$

② Suppose $f: \mathbb{R}^n \rightarrow \mathbb{R}$

$$g: \mathbb{R}^m \rightarrow \mathbb{R}$$

$$A \in \mathbb{R}^{n \times m}$$

$$b \in \mathbb{R}^m$$

$$g(x) = f(Ax+b)$$

$$\nabla^2 g(x) = A^T \nabla^2 f(Ax+b) A$$

③ Define $\tilde{f}(t) = f(x+tv)$

$$\nabla^2 \tilde{f}(t) = \tilde{f}''(t) = v^T \nabla^2 f(x+tv)v$$

$$\text{For } t=0: \nabla^2 \tilde{f}(0) = v^T \nabla^2 f(x)v$$

Example: $f(x) = \frac{1}{2} x^T P x + q^T x + r$

$$P \in \mathbb{S}^n, q \in \mathbb{R}^n, r \in \mathbb{R}$$

$$\nabla f(x) = \underline{Px+q}$$

$$\nabla^2 f(x) = D(\nabla f(x)) = P$$

assume: $f(x) = \frac{1}{2} \boxed{ax^2} + bx + c$

$$\underline{f''(x) = a}$$

Newton's Method

It is "supposed" to be a descent method with iterations given by:

$$x^{(k+1)} = x^{(k)} + t^{(k)} \Delta x_{nt}$$

where: $\Delta x_{nt} = -[\nabla^2 f(x)]^{-1} \nabla f(x)$ Newton direction



Based on this, it requires:

① Function f has to be twice differentiable

↳ Typically we assume $f \in C^2$

↳ Twice continuously differentiable

② $\nabla^2 f(x)$ must be invertible $\Rightarrow \text{rank}(\nabla^2 f(x)) = n$

↳ Typical requirement is that it is invertible

over every $x \in \mathbb{R}^n$

③ In order for Newton's method to be a descent method, we require that

$$\nabla^2 f(x) > 0 \quad (\text{Positive definite})$$

(Remember: $A > 0 \iff A^{-1} > 0$)

There are two ways to handle ③

① Assume $\nabla^2 f(x) > 0 \quad \forall x \in \mathbb{R}^n$

↳ strongly convex functions

② What about non convex functions?



In that case, we first run gradient descent for a number of iterations till $\|\nabla f(x)\|_2$ is small and then we switch to Newton iterations.

Even when $\nabla^2 f(x) > 0 \forall x \in \mathbb{R}^n$, Newton's has some drawbacks:

- ① Compute and store $\nabla^2 f(x)$
- ② Compute inverse of $[\nabla^2 f(x)]^{-1}$

so Why use it? It is extremely fast in the right regions (to be shown later).

→ Ways to deal with these issues

- ① Quasi-Newton method
- ② Approximate the Hessian by looking/exploring the structure of the problem (in a fast way).

$$\text{E.g., } \nabla f(x) = \begin{bmatrix} f_1(x_1) \\ f_2(x_2) \\ \vdots \\ f_n(x_n) \end{bmatrix}$$

$$\Rightarrow \nabla^2 f(x) = \begin{bmatrix} f_{11}(x_1) & f_{12}(x_2) & \cdots & f_{1n}(x_n) \\ f_{21}(x_1) & f_{22}(x_2) & \cdots & f_{2n}(x_n) \\ \vdots & \vdots & \ddots & \vdots \\ f_{n1}(x_1) & f_{n2}(x_2) & \cdots & f_{nn}(x_n) \end{bmatrix}$$

Interpretations of Newton's Method

- ① Minimizer of the second-order approximation of f at x

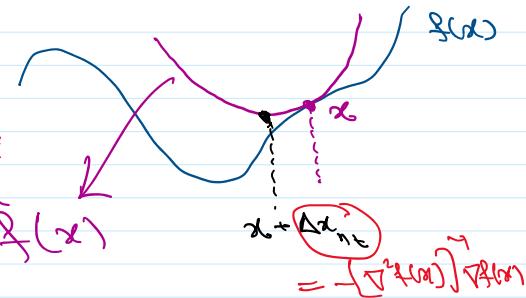
① Minimizer of the second-order approximation of f at x

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

write $y = x + v$

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

$$\begin{aligned} \text{argmin}_v &= -\left[\nabla^2 f(x)\right]^{-1} \nabla f(x) \\ &= -\left[\nabla^2 f(x)\right]^{-1} \nabla f(x) \end{aligned}$$



Compute ∇_v and set it equal to 0.

$$\begin{aligned} \nabla_v &\left(f(x) + \nabla f(x)^T v + \underbrace{\frac{1}{2} v^T \nabla^2 f(x) v}_{\nabla^2 f(x)v} \right) \\ &= 0 + \nabla f(x)^T + \nabla^2 f(x)v = 0 \end{aligned}$$

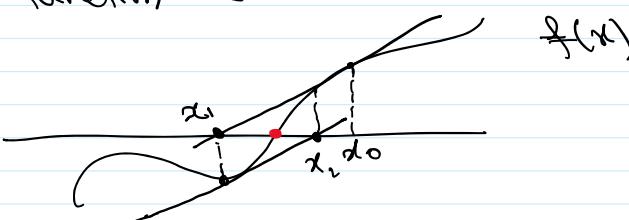
$$\begin{aligned} \nabla^2 f(x)v &= -\nabla f(x) \\ v^* &= -\left[\nabla^2 f(x)\right]^{-1} \nabla f(x) \end{aligned}$$

② Newton's method is also tied to the idea of approximating the gradient $\nabla f(x)$ by a linear function and then finding the root of that linear function.

Stationary point of a function is

when $\nabla f(x) = 0$

Given $f(x) : \mathbb{R} \rightarrow \mathbb{R}$



$$\underline{f(y)} \approx f(x) + \dot{f}(x)(y-x)$$

$$\nabla f(y) \sim \nabla f(x) + \nabla^2 f(x)(y-x)$$

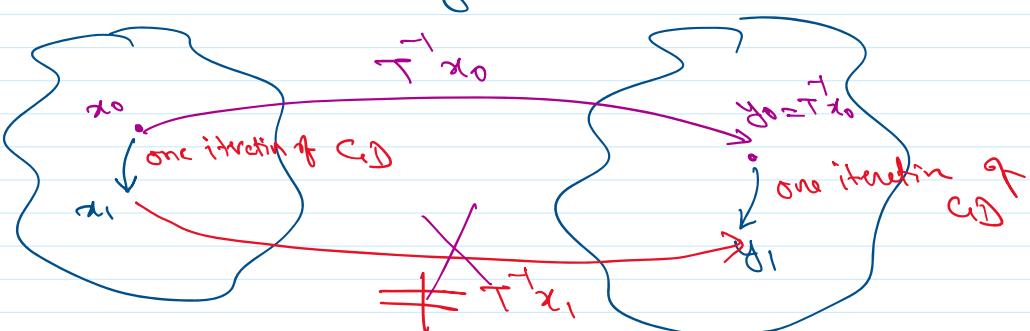
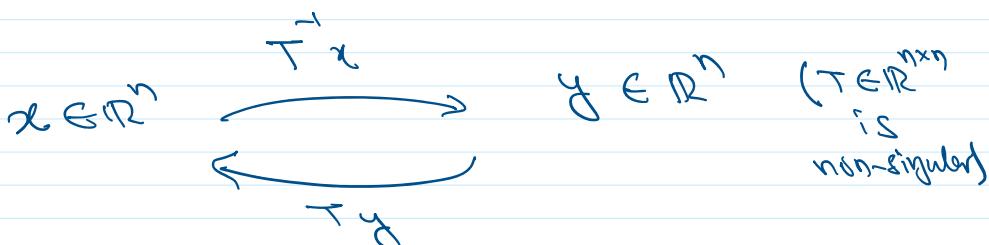
$$\nabla \mathcal{L}(y) \approx \nabla \mathcal{L}(x) + \nabla^2 \mathcal{L}(x)(y-x)$$

Put $y = x + v$

$$\underbrace{\nabla \mathcal{L}(x+v)}_{=0} = \nabla \mathcal{L}(x) + \nabla^2 \mathcal{L}(x)v$$

$$\nabla^2 \mathcal{L}(x)v = -\nabla \mathcal{L}(x)$$

$$v = -[\nabla^2 \mathcal{L}(x)]^{-1} \nabla \mathcal{L}(x)$$



Gradient descent, in general, is not affine invariant. Coordinate system in gradient descent affects the algorithmic performance.

Affine invariance of Newton's Step

Suppose $T \in \mathbb{R}^{n \times n}$ is non-singular and

let $y = T^{-1}x$; $x = Ty$

$$\mathcal{L}(x) : \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\bar{\mathcal{L}}(y) = \mathcal{L}(Ty) : \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\bar{f}(y) = f(Ty) : \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\text{Then: } x + \Delta_{\text{int}} = T(y + \Delta y_{\text{int}})$$

Basic Assumption

Either we are close to a local optimum or we are working with the case $\nabla^2 f(x) > 0 \forall x \in \mathbb{R}^n$.

Newton Decrement

$$\gamma(x) = \left[\nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x) \right]^{1/2}$$

is called Newton decrement.

- ① Used in analysis
- ② Used in stopping criterion of Newton's method

- $\gamma(x)$ is a scalar
- $\gamma(x) > 0$ since $\nabla^2 f(x)^{-1} > 0$

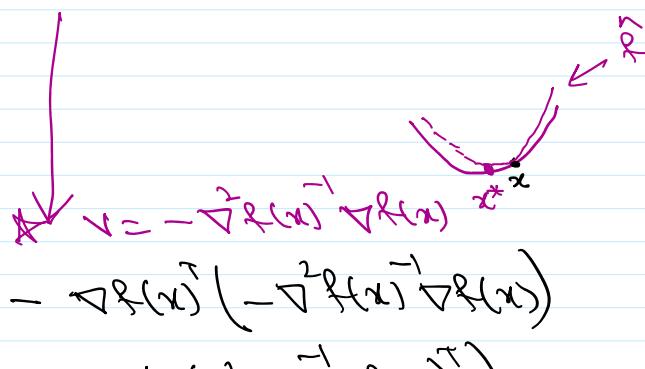
$\gamma(x)$ allows us to approximate how close we are to a local minimum (x^*).

$$\text{Recall: } \hat{f}(x+v) = f(x) + \nabla f(x)^T v + \frac{1}{2} v^T \nabla^2 f(x) v$$

$$\underbrace{f(x)}_{\text{function at } x} - \underbrace{\hat{f}(x+v)}_{\text{Quadratic approximation at } x} = -\nabla f(x)^T v - \frac{1}{2} v^T \nabla^2 f(x) v$$

function
at x

Quadratic
approximation
at x

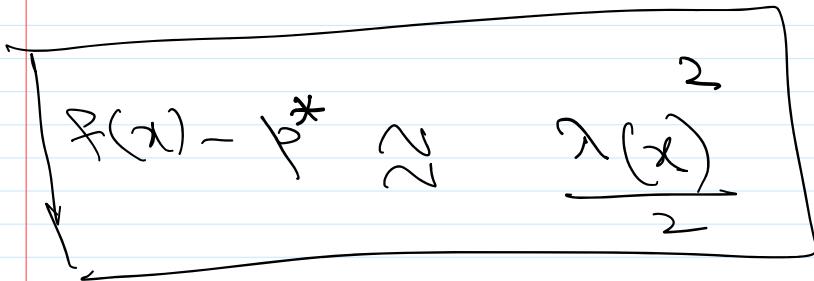


$$f(x) - \min_v \hat{f}(x+v) \approx -\nabla f(x)^T \left(-\nabla^2 f(x)^{-1} \nabla f(x) \right)$$

$$f(x) = \min_v f(x+v) \approx -\nabla f(x)^T (-\nabla^2 f(x) \nabla f(x))$$

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

$$\nabla^2 f(x) (\nabla^2 f(x)^{-1} \nabla f(x))$$



When our current iteration x is very close to x^* , $\frac{\|\nabla f(x)\|^2}{2}$ gives us an estimate of how far we are from the local minimum.

Newton Decrement: $\lambda(x) = \left\{ \nabla^2 f(x)^{-1} \nabla f(x) \right\}_2^{1/2}$

When we are close to a minimizer p^*

$$\frac{\lambda(x)}{2} \Rightarrow \text{Estimate of how close we are to } p^*$$

$$f(x) - f^* \approx \frac{\lambda(x)}{2}$$

Some facts: ① $\nabla f(x)^T \Delta x_{nt} = -\lambda(x)^2$

$$-\nabla^2 f(x)^{-1} \nabla f(x)$$

$$\|\nabla f(x)\|_2 \|\Delta x_{nt}\|_2 \cos(\theta)$$

② $\lambda(x) = (\Delta x_{nt}^T \nabla^2 f(x) \Delta x_{nt})^{1/2}$

P-Quadratic norm of a vector when $P \in S_{++}^n$

$$\|x\|_P = (x^T P x)^{1/2}$$

$\rightarrow \lambda(x) = \|\Delta x_{nt}\|_{\nabla^2 f(x)}$

③ $\lambda(x)$ is invariant under affine transformations of the function $f(x)$

Newton's Step:

$$x^{(u+1)} = x^{(u)} + t^{(u)} \Delta x_{nt}$$

$\underbrace{-\nabla^2 f(x)^{-1} \nabla f(x)}$
 $\underbrace{\nabla^2 f(x) > 0}$

when $t^{(u)} = 1 \Rightarrow$ Full Newton Step.

↳ Full Newton Step should be taken only when x is very close to the local (or global) minimum. Otherwise, the method can diverge.

Newton's Method (Assuming we are close to a local minimum or are working with a convex problem)

Input: $x^{(u)} \in \text{dom} f$; tolerance $\epsilon > 0$

Initialize: $k \leftarrow 0$

Repeat:

$$\textcircled{1} \quad \Delta x_{nt}^{(u)} \leftarrow -\nabla^2 f(x^{(u)})^{-1} \nabla f(x^{(u)}) \quad // \text{Compute Newton's direction}$$

↳ make sure
 $\nabla^2 f(x^{(u)}) > 0$

$$\textcircled{2} \quad \text{Compute } \lambda^2(x^{(u)}) \leftarrow \nabla F(x^{(u)})^\top \nabla^2 f(x^{(u)})^{-1} \nabla F(x^{(u)}) \quad // \text{Compute Newton decrement square}$$

$$\textcircled{3} \quad \text{If } \frac{\lambda^2(x^{(u)})}{\lambda^2(x^{(u-1)})} \leq \epsilon$$

2
break

④ Line Search: choose $t^{(k)}$ using back tracking
line search

⑤ Update: $x^{(k+1)} \leftarrow x^{(k)} + t^{(k)} \Delta x_{nt}^{(k)}$

⑥ $k \leftarrow k + 1$

~~~~~ X ~~~~~ X ~~~~~

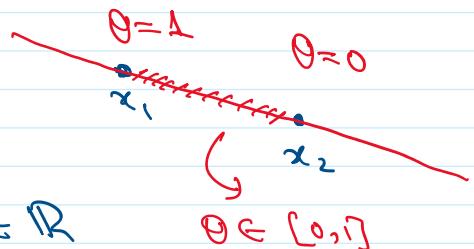
Chapter 2: Convex Sets

① Line: Given  $x_1, x_2 \in \mathbb{R}^n$

$$y = \theta x_1 + (1-\theta)x_2; \theta \in \mathbb{R}$$

$$\downarrow y(\theta)$$

↳ one-dimensional linear function



$y(\theta)$  is a line passing through  $x_1$  and  $x_2$

② Line Segment: when  $\theta \in [0,1]$ ,

$y(\theta) = \theta x_1 + (1-\theta)x_2$ ; is the

line segment } joining  $x_1$  and  $x_2$

... . o / ~ ~ ~

$$\downarrow$$

$$= x_2 + \theta(x_1 - x_2)$$

$\Theta \Rightarrow$  tells us how far we are from  
 $x_2$  in the direction  $(x_1 - x_2)$

③ Affine Set: A set  $C \subseteq \mathbb{R}^n$  is affine if the line through any distinct points in  $C$  lies in the set  $C$ .

$\forall x_1, x_2 \in C$  and  $\theta \in \mathbb{R}$

$$\theta x_1 + (1-\theta)x_2 \in C$$

$\hookrightarrow$  Linear combinations of points in  $C$  lie in  $C$  provided the coefficients of the linear combination sum to 1.

$\Downarrow$  simple induction argument

If  $C$  is an affine set,  $x_1, \dots, x_k \in C$   
 $\theta_1, \dots, \theta_k \in \mathbb{R}$  s.t.  $\sum_{j=1}^k \theta_j = 1$  then

$$\theta_1 x_1 + \theta_2 x_2 + \dots + \theta_k x_k \in C$$

### Affine Combination of Points

Given  $x_1, \dots, x_n \in \mathbb{R}^n$  and  $\theta_1, \dots, \theta_n \in \mathbb{R}$   
s.t.  $\sum_{j=1}^n \theta_j = 1$

$\theta_1x_1 + \theta_2x_2 + \dots + \theta_kx_k$  is called  
an affine combination.

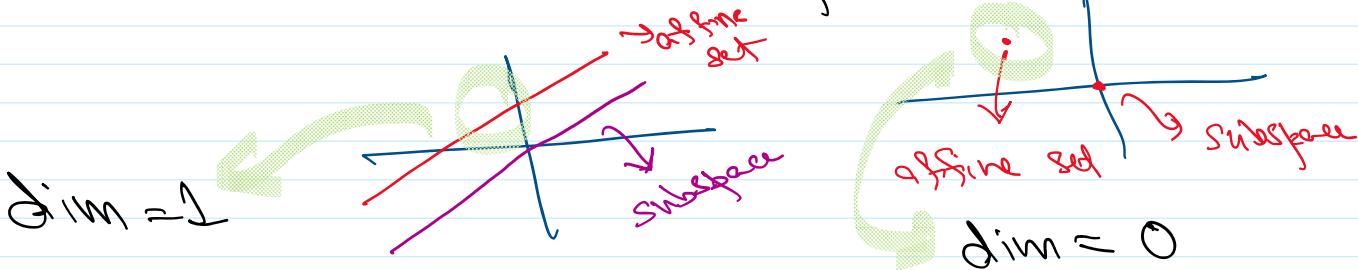
### Connection to a Subspace

An affine set is a subspace shifted from the origin.

If  $C$  is an affine set and  $x_0 \in C$  for any  $x_0$  then

$$V = C - x_0 = \{x - x_0 : x \in C\}$$

Then  $V$  is a subspace.



### Dimension of an Affine Set

It is the dimension of the subspace associated with  $C$ .

Ex:  $C = \{x : Ax = b ; b \in \mathbb{R}^n\}$

Is  $C$  affine?

Approach 1: Use the definition of an affine set.

Let  $x_1$  and  $x_2 \in C$

Show that  $\theta x_1 + (1-\theta)x_2 \in C$

$$x_1 \in C \Rightarrow Ax_1 = b$$

$$x_2 \in C \Rightarrow Ax_2 = b$$

$$\theta x_1 + (1-\theta)x_2 \Rightarrow \in C ?$$

Show  $\Rightarrow A(\theta x_1 + (1-\theta)x_2) = b$

$$= \underbrace{\theta Ax_1}_b + \underbrace{(1-\theta)Ax_2}_b$$

$$= \theta b + (1-\theta)b = b \checkmark$$

Approach 2:-  $\{x: Ax = 0\}$  is a subspace  
null space of A

$\{x: Ax = b\}$  is simply a shifted version of  
 $\{x: Ax = 0\} \Rightarrow C$  is an affine set.

Dimensionality of C = dimension of the null space of A.

Affine Hull of points in a set in  $\mathbb{R}^n$

Given a set C, which may not be affine,  
the affine hull of C is defined as

Given a set  $C$ , which may not be affine,  
the affine hull of  $C$  is defined as:

$$\text{aff } C = \left\{ \theta_1 x_1 + \dots + \theta_n x_n : x_1, \dots, x_n \in C \right. \\ \left. \theta_1 + \theta_2 + \dots + \theta_n = 1 \right\}$$

↪ Affine hull of a set  $C$  is the smallest affine set that contains  $C$ . That is, let  $S$  be any affine set such that:

$$C \subseteq S. \text{ Then.}$$

$$\text{aff } C \subseteq S$$

Examples :

$$\text{aff } \left( \begin{array}{c} \cdot \\ \cdot \\ \hline \end{array} \right) = \begin{array}{c} \cdot \\ \cdot \\ \hline \end{array}$$

$$\text{aff } \left( \begin{array}{c} \cdot \\ \cdot \\ \hline \end{array} \right) = \begin{array}{c} \cdot \\ \cdot \\ \hline \end{array}$$

$$\text{aff } \left( \begin{array}{cc} \cdot & \cdot \\ \cdot & \cdot \\ \hline \end{array} \right) = \begin{array}{c} \cdot \\ \cdot \\ \hline \end{array}$$

$$\text{aff } \left( \begin{array}{c} \cdot \\ \cdot \\ \hline \end{array} \right) = \mathbb{R}^2$$

Reading: 2.1.1 and 2.1.2 (BV)



### Convex Set

A set  $C$  is convex if the line segment between any two points in  $C$ , lies in  $C$ .

$$\text{if } x_1, x_2 \in C \text{ and } \theta \in [0,1]$$

$$\theta x_1 + (1-\theta)x_2 \in C$$

- Every affine set is a convex set



$\Rightarrow$  Convex



$\Rightarrow$  nonConvex



$\Rightarrow$  Convex



$\Rightarrow$  nonConvex

nonConvex

Linear combination

Subspaces

$$\sum_{i=1}^k \theta_i = 1$$

$$\begin{array}{c} \sum_{i=1}^k \theta_i = 1 \\ \theta_i \geq 0 \end{array}$$

$\subset$  affine sets  $\subset$  convex sets

### Convex Combination of Points

Given  $x_1, x_2, \dots, x_n \in C$  and  $\theta_1, \dots, \theta_n \geq 0$

$$\text{s.t. } \sum_{j=1}^k \theta_j = 1$$

$\theta_0x_0 + \dots + \theta_kx_k$  is called a convex combination of the points in  $C$ .

A convex set implies that any convex combination of points in the set lie in the same set.

### Convex hull

The convex hull of a set  $C$  is the set of all convex combinations of points in  $C$ .

$$\text{conv } C = \left\{ \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n : x_i \in C; \theta_i \geq 0; \sum_{i=1}^n \theta_i = 1 \right\}$$

Convex hull of  $C$  is convex  $\Leftrightarrow$  It is the smallest convex set that contains  $C$ .

### Examples :

$$\text{Conv} \left( \begin{array}{c} \cdot \\ - \end{array} \right) = \begin{array}{c} \cdot \\ - \end{array}$$

$$\text{Conv} \left( \begin{array}{c} \cdot \\ \cdot \\ - \end{array} \right) = \begin{array}{c} \cdot \\ \nearrow \\ - \end{array}$$

$$\text{Conv} \left( \begin{array}{c} \cdot \\ \cdot \\ \cdot \\ - \end{array} \right) = \begin{array}{c} \nearrow \\ - \end{array}$$

$$\text{Conv} \left( \begin{array}{c} \cdot \\ \cdot \\ \cdot \\ \bullet \\ - \end{array} \right) = \begin{array}{c} \bullet \\ - \end{array}$$

$$\text{Conv} \quad \left( \begin{array}{c} \alpha \\ \beta \end{array} \right) = \begin{array}{c} \alpha \\ \beta \end{array}$$

Convex combination  $\Rightarrow$

$$\left. \begin{array}{l} \theta_i \geq 0, i=1\dots,k \\ \sum_{i=1}^k \theta_i = 1 \end{array} \right\} \text{Probability mass function}$$

$$\left. \begin{array}{l} p_x(x) \geq 0 \\ \int p_x(x) dx = 1 \end{array} \right\} \text{Probability density functions}$$

$$E[X] = \sum_{i=1}^k p_x(x_i) x_i \geq 0$$

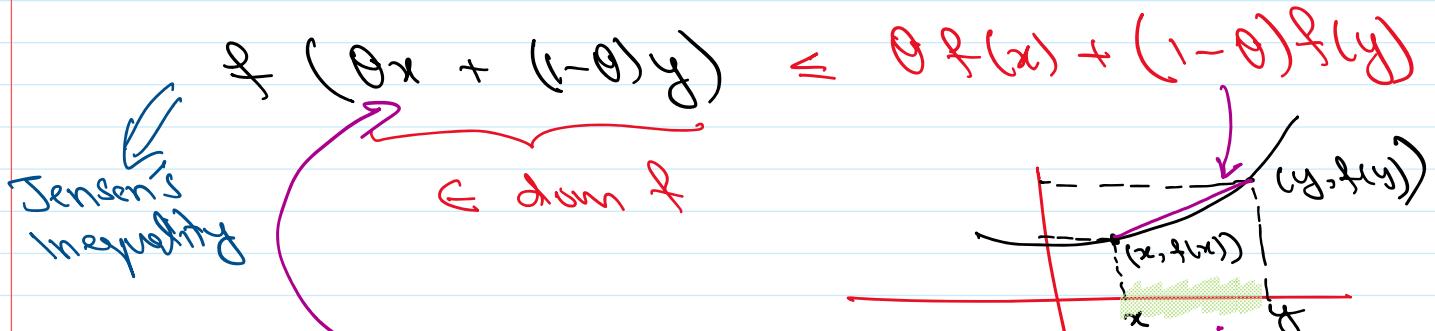
Reading: 2.1.4 (BV)

## Convex functions

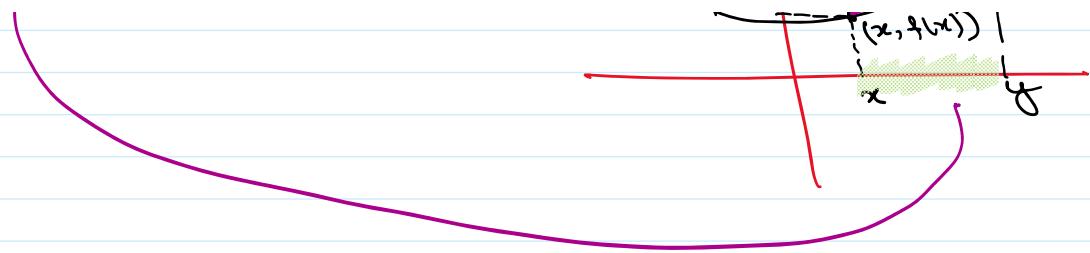
Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  with  $\text{dom } f$  being convex.

Then  $f$  is termed a convex function if

$\forall x, y \in \text{dom } f$  and  $\theta \in [0, 1]$



Impression



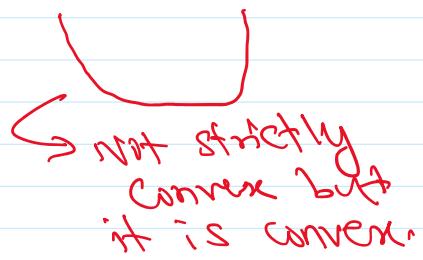
The chord (line segment) connecting  $(x, f(x))$  and  $(y, f(y))$  should lie above the function between  $x$  and  $y$ .

### Strictly Convex Function

If  $f(\theta x + (1-\theta)y) < \theta f(x) + (1-\theta)f(y)$

&  $x, y \in \text{dom } f$  then  $f$  is called a strictly convex function.

Ex: A linear function is convex but it is not strictly convex.



Concave functions: If  $-f$  is

convex then  $f$  is called concave  
(similarly Strictly Concave).

What if  $\text{dom } f \neq \mathbb{R}^n$ ?

Extensions of convex functions on all  $\mathbb{R}^n$

Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be convex on  $\text{dom } f$ .

The Convex extension or  $\tilde{f}$  on  $\mathbb{R}^n$  is defined

The Convex extension of  $f$  on  $\mathbb{R}^n$  is defined as:

$$\tilde{f}(x) = \begin{cases} f(x), & x \in \text{dom } f \\ \infty, & x \notin \text{dom } f \end{cases}$$

↳ using this trick, we can ignore that  $\text{dom } f \neq \mathbb{R}^n$

$$\text{dom } f = \{x : \tilde{f}(x) < \infty\}$$

Ex:  $f(x) = -\log x$

$\text{dom } f = (0, \infty)$



### Constrained Optimization trick

Let  $C$  be a convex set

Suppose we need to solve

$$\min_{x \in C} f(x); \quad f: \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\tilde{I}_C(x) = \begin{cases} 0, & x \in C \\ \infty, & x \notin C \end{cases} \Rightarrow \text{Convex function}$$

$$\min_{x \in \mathbb{R}^n} [f(x) + \tilde{I}_C(x)]$$

Reading: 3.1.1 and 3.1.2 (BV)

Convex Function, Convex Combination, and Probability

## Convex Function, Convex Combination, and Probability

If  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is convex

$$\Rightarrow f(\theta_1 x_1 + \dots + \theta_n x_n) \leq \theta_1 f(x_1) + \dots + \theta_n f(x_n)$$

$\theta_i \geq 0, \sum_i \theta_i = 1 \quad \hookrightarrow E[X] \quad E[f(X)]$

Let  $x_i$ 's be the values that a random variable  $X$  takes and  $\theta_i$  are the probabilities

$$P(X=x_i) = \theta_i$$

If  $f$  is convex  $\Rightarrow f(E[X]) \leq E[f(X)]$   
 $\hookrightarrow$  Jensen's Inequality

$\hookrightarrow$  This holds even when  $X$  is continuous random variable.

Reading: 3.1.8 (BV)

## Equivalent Characterizations of Convex Functions

① A function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is convex if and only if its restriction to any line in  $\mathbb{R}^n$  is convex:  
domf is convex and

Define:  $g(t) : \mathbb{R} \rightarrow \mathbb{R}$

$$g(t) = f(x + tv) \quad \forall x + tv \in \text{dom}f$$

Then  $f$  is convex  $\Leftrightarrow g(t)$  is convex  
for every  $x$  and  $v$ .

## ② First-order Condition of Convexity

Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be differentiable on  $\text{dom } f$  and  $\text{dom } f$  is open.

$f: \mathbb{R}^n \rightarrow \mathbb{R}$  is convex if and only if  $\text{dom } f$  is convex and

$$f(y) \geq f(x) + \nabla f(x)^T (y-x) \quad \text{--- (2)}$$

$\forall x, y \in \text{dom } f$

first-order approximation  
(linear approximation)  
of  $f$  around

↳ The first-order approximation must be a uniform underestimator of  $f$ .

## Global optimality condition for convex functions

Let  $x_0 \in \text{dom } f$  be such that  $\nabla f(x_0) = 0$

then  $x_0$  is a global minimizer of  $f$ .

i.e.

$$f(x_0) \leq f(y) \quad \forall y \in \text{dom } f$$

Proof: Take  $x = x_0$  in (2) (first-order convexity condition)

$$f(y) \geq f(x_0) + \nabla f(x_0)^T(y - x_0)$$

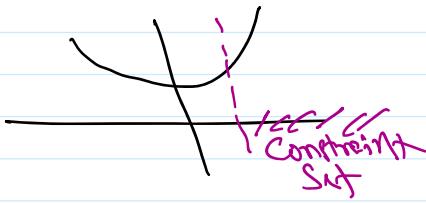
$$\Rightarrow f(x_0) \leq f(y) \quad \forall y \in \text{dom} f$$

~~Contradiction~~

### Unconstrained optimization

$x_0$  is a global minimizer of  $f$  if and only if

$$\nabla f(x_0) = 0$$



### Strictly convex functions

$$f(y) > f(x) + \nabla f(x)^T(y - x)$$

A strictly convex function can only have a unique minimizer.

Indeed: Let  $x_1$  and  $x_2$  be two global minimizers with  $\nabla f(x_1) = 0 = \nabla f(x_2)$

$$f(x_1) < f(y) \quad \forall y \in \text{dom} f$$

$$f(x_1) < f(x_2)$$

$\Rightarrow$  Contradiction  $\square$

### ③ Monotonicity of gradients

Suppose  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is differentiable with  $\text{dom} f$  being convex. Then  $f$  is convex if and only if

$\rightarrow$

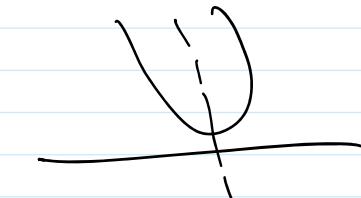
convex. Then  $f$  is convex if and only if

$$(\nabla f(x) - \nabla f(y))^\top (y-x) \geq 0 \quad \forall x, y \in \text{dom}f$$

$> 0 \Rightarrow \text{strict convexity}$

This generalizes the concept of monotonicity of functions to  $\mathbb{R}^n$

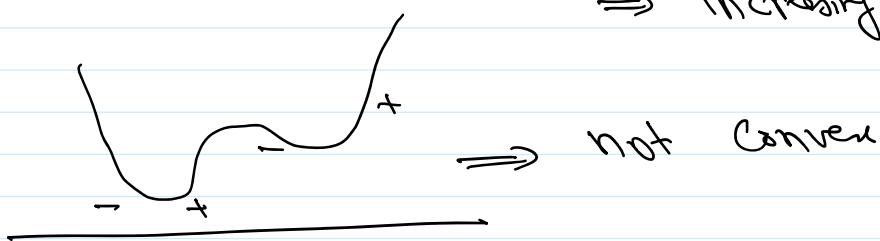
E.g.,



$$f(x) = x^2$$

$$\Rightarrow f'(x) = 2x$$

$\Rightarrow$  Increasing  $\Rightarrow$  convex



#### ④ Second-order Condition of Convexity

Suppose  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is twice differentiable with  $\text{dom}f$  being open.

$f$  is convex if and only if  $\text{dom}f$  is convex and

$$\underbrace{\nabla^2 f(x)}_{\text{Positive semi definite.}} \geq 0 \quad \forall x \in \text{dom}f$$

Positive semi definite.

Basically, the function at every point  $x$  has non-negative curvature.

Concave:  $\nabla^2 f(x) \leq 0$

Strictly Convex functions  $\nabla^2 f(x) > 0$

Ex.  $f(x) = x^4 \Rightarrow$  Strictly Convex.  
but  $f''(x) = 0$  at  $x=0$

## Quadratic Functions

$f: \mathbb{R}^n \rightarrow \mathbb{R}$  with  $\text{dom } f = \mathbb{R}^n$

$$f(x) = \frac{1}{2} x^T P x + q^T x + r$$

$P \in \mathbb{R}^{n \times n}$

↳ Quadratic functions

$$\nabla^2 f(x) = P$$

A quadratic function is convex if and only if

$$P \succeq 0 \quad (\text{P is S}^n_+)$$

$f(x)$  is concave when  $P \preceq 0$  (negative semi-definite)

when  $P$  is neither  $\succeq 0$  or  $\preceq 0 \Rightarrow$  It is neither concave nor convex

1D:  $f(x) = ax^2$ ;  $f(x) = x^2$

↳ convex

$$f(x) = -x^2 \Rightarrow \text{concave}$$

2D:  $f(x) = x_1^2 - x_2^2$

## Unconstrained Minimization of Convex Functions

$$\min_x f(x) \quad \text{s.t. } f: \mathbb{R}^n \rightarrow \mathbb{R}$$

If  $f$  is convex  $\Leftrightarrow x^* : \nabla f(x^*) = 0$

Necessary and sufficient condition for solution of a convex optimization problem corresponds to solving the equation  $\nabla f(x) = 0$

e.g.,  $f(x) = (x-4)^2 \Rightarrow f'(x) = 2(x-4) = 0$   
 $\Rightarrow x^* = 4$

Why then we need optimization algorithms?

① Finding a solution to  $\nabla f(x) = 0$  is not always possible analytically.

e.g.,  $\min_x f(x) = \log \left( \sum_{i=1}^m \exp(a_i^\top x + b_i) \right)$   
↳ geometric program

$$\nabla f(x) = \frac{1}{\sum_{i=1}^m \exp(a_i^\top x + b_i)} \sum_{i=1}^m \exp(a_i^\top x + b_i) a_i$$

② Even when an analytical solution exists, computing it numerically in a direct fashion might be too costly.

e.g.,  $\min_x f(x) = \|Ax - b\|_2^2$

$$\nabla f(x) = 2A^\top(Ax - b)$$

$$\Rightarrow \nabla f(x) = 0 \Leftrightarrow \underbrace{A^T A x = A^T b}_{\text{Let's even assume } (A^T A) \text{ is full rank.}}$$

$$x^* = (A^T A)^{-1} A^T b$$

## Convex Optimization Problems with and without Solution

$$\textcircled{1} \quad f(x) = a^T x$$

$\min_x f(x) \Rightarrow$  no solution  $\Rightarrow$  function is unbounded below

$$\textcircled{2} \quad f(x) = e^{-x}$$

$\min_x f(x) \Rightarrow$  does not exist

$\hookrightarrow$  function is bounded below but min does not exist

$$\inf_x \{f(x)\} = e^{-\infty} = 0$$

Supremum (sup) and infimum (inf)  $\Rightarrow$  Generalizations of max and min of sets

Let  $C \subseteq \mathbb{R}$ .

We say that a scalar value ' $u$ ' is an upperbound on the set  $C$  if:

$$\forall x \in C : x \leq u$$

① If the set of upper bounds is empty  $\Rightarrow C$  is unbounded from above

$C = \{1, 2, 3\} \Rightarrow C$  is unbounded from above.

② If  $C = \emptyset \Rightarrow$  set of upper bounds is  $\mathbb{R}$

The set of upper bounds is always of the form  $[b, \infty)$

e.g.,  $C = \{1, 2, 3\} \Rightarrow$  set of upper bounds =  $[3, \infty)$

$\sup(C)$  is defined as the least upper bound in the set of upper bounds.

①  $\sup \{1, 2, 3\} = 3 = \max$

$\sup = \max$  when  $\sup$  belongs to the set  $C$

②  $\sup \left\{ 1 - \frac{1}{n} \right\}_{n=1}^{\infty} = 1$   
max  $\Rightarrow$  does not exist

\*  $\sup(\emptyset) = -\infty$

what if  $C$  is unbounded?

\*  $\sup(C) = \infty$

Infinum of a set  $\Rightarrow$  deals with the case of lower bounds on sets

The set of lower bounds takes the form

$(-\infty, l]$

$\inf(C) \Rightarrow$  The greatest lower bound on the set

c.

e.g.)  $C = \left\{ 1 + \frac{1}{n} \right\}_{n=1}^{\infty}$

$$\inf(C) = 1 \neq \min(C)$$

↳ does not exist

$$\inf \{1, 2, 3\} = 1 = \min(C)$$

\*  $\inf(\emptyset) = \infty$

\*  $\inf(C) = -\infty \Rightarrow$  when  $C$  is unbounded from below

$$\inf(C) = -\sup(-C)$$

Optimization problems should technically be initially written as inf problems, till we are convinced that  $\inf = \min$

$$\inf_x a^T x = -\infty \quad \text{if } a \neq 0$$

$$f^* = \inf_x e^{-x} = 0 \neq \boxed{\min_x e^{-x}}$$

↳ This is not solvable

$$\arg \inf_x e^{-x} = \infty \text{ or no solution}$$

↳ There is no  $x^*$  such that  $f(x^*) = f^*$

Example of a Quadratic

$$f(x) = \frac{1}{2} x^T P x + q^T x + r ; \quad P \succcurlyeq 0$$

$\Rightarrow f(x)$  is convex

$$\nabla f(x) = Px + q$$

If  $\min_x f(x)$  has a solution

$$\Leftrightarrow \nabla f(x) = 0$$

$$Px + q = 0 \Leftrightarrow Px = -q$$

The Solution Set is  $\{x : Px = -q\}$

What if this is  
not the case?

$-q$  is in the column  
space of  $P$

$\Rightarrow f(x)$  is unbounded from  
below and there is no solution.

If  $P$  has full rank  $\Rightarrow$  The Solution set  
is non-empty; in fact there is a unique solution.  
 $x^* = -P^{-1}q$

$$P \succ 0 \Rightarrow f(x) \text{ is strictly convex}$$

Strongly Convex Functions

SC functions  $\subset$  Strictly Convex  $\subset$  Convex functions

Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  with domf being convex.

$f$  is strongly convex with parameter  $m > 0$  if:

$$\Leftrightarrow \textcircled{1} \quad f(\theta x + (1-\theta)y) \leq \theta f(x) + (1-\theta)f(y)$$

$$-\frac{\theta(1-\theta)}{2} \|x-y\|_2^2$$

A function is strongly convex if and only if it grows at least as fast as a quadratic function with growth parameter 'm'.

From \textcircled{1}, we see that

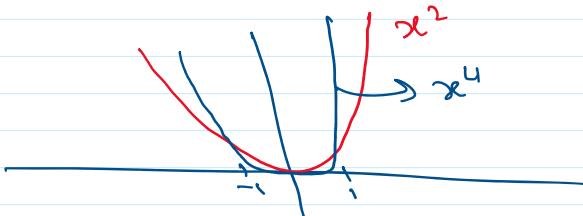
$$f(\theta x + (1-\theta)y) < \theta f(x) + (1-\theta)f(y)$$

$\Rightarrow$  SC functions are strictly convex also

$\Rightarrow$  SC functions have a unique minimizer.

E.g.  $f(x) = x^2 \rightarrow$  strongly convex

$f(x) = x^4 \rightarrow$  strictly convex



$\Rightarrow$  But not strongly convex

$\Leftrightarrow \textcircled{2}$  Second-order definition of strong convexity

Suppose  $f \in C^2$  (twice continuously differentiable) and dom $f$  is open

$f(x)$  is SC  $\Leftrightarrow \nabla^2 f(x) \geq mI \quad \forall x \in \text{dom}f$

$\nabla^2 f(x) \geq mI \Leftrightarrow \nabla^2 f(x) - mI \geq 0$

$\Leftrightarrow \lambda_{\min}(\nabla^2 f(x)) \geq m > 0$

e.g.,  $f(x) = x^4 \Rightarrow f'(x) = 4x^3$   
 $f''(x) = 12x^2$

$f''(x) > 0 \quad \forall x \neq 0$  but

$f''(x) = 0$  for  $x=0$ .

Even for  $x \neq 0$   $f''(x) \not> m$  for some  $m > 0$ .