

Session 9: Convex functions

Optimization and Computational Linear Algebra for Data Science

Contents

1. Recap of the videos
2. Convex sets and convex functions
3. Convex functions and derivatives
4. Jensen's inequality

Optimization

In machine learning, we often have to minimize functions

$$f(\theta) = \text{Loss}(\text{data}, \text{model}_\theta) \quad \text{with respect to } \theta \in \mathbb{R}^n.$$

model's parameters

- ❖ For $n = 1, 2$, one could plot f to find the minimizer.
- ❖ This is intractable for larger dimension.

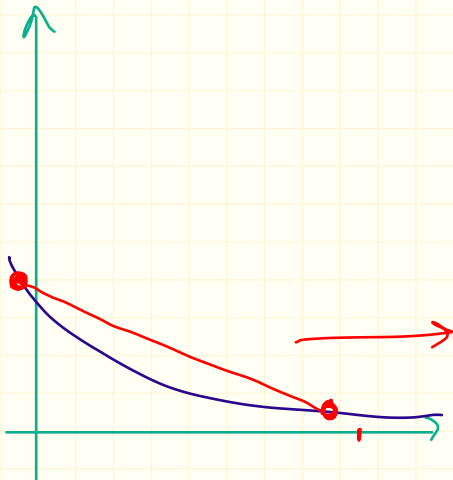
We will

- ❖ focus on **convex cost functions** f .
- ❖ study gradient descent algorithms to minimize f .

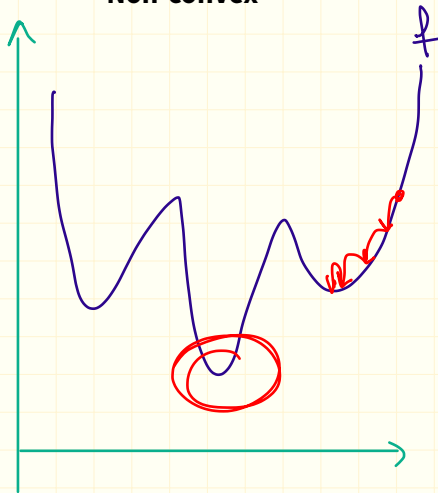
Convex vs non-convex

Convex

e^{-x}



Non-convex

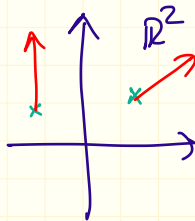


Gradient/Hessian

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

▣ Gradient at $x \in \mathbb{R}^n$:

$$\nabla f(x) = \begin{pmatrix} \frac{\partial f}{\partial x_1}(x) \\ \vdots \\ \frac{\partial f}{\partial x_n}(x) \end{pmatrix} \in \mathbb{R}^n$$



▣ Hessian at $x \in \mathbb{R}^n$:

$$H_f(x) = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2}(x) & \frac{\partial^2 f}{\partial x_1 \partial x_2}(x) & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n}(x) \\ \frac{\partial^2 f}{\partial x_2 \partial x_1}(x) & \frac{\partial^2 f}{\partial x_2^2}(x) & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n}(x) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1}(x) & \frac{\partial^2 f}{\partial x_n \partial x_2}(x) & \cdots & \frac{\partial^2 f}{\partial x_n^2}(x) \end{pmatrix} \in \mathbb{R}^{n \times n}$$

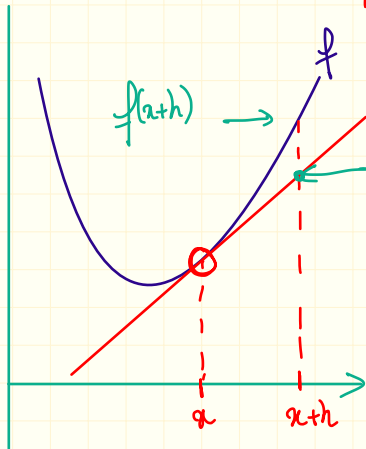
Taylor's formulas

Let $x \in \mathbb{R}^n$. Heuristically, for $h \in \mathbb{R}^n$ "small", we have

$$f(x+h) \simeq f(x) + \langle \nabla f(x), h \rangle$$

$u \cdot v$
 $\langle u, v \rangle =$ dot product between u and v

not mathematically correct



"Rigorous Taylor"

$$f(x+h) = f(x) + \nabla f(x) \cdot h + \text{Err}_1(h)$$

$$\frac{\text{Err}_1(h)}{\|h\|} \xrightarrow{h \rightarrow 0} 0$$

0,001

Taylor's formulas

Let $x \in \mathbb{R}^n$. Heuristically, for $h \in \mathbb{R}^n$ "small", we have

$$f(x+h) \simeq f(x) + \langle \nabla f(x), h \rangle + \frac{1}{2} h^\top H_f(x) h.$$

$$+ \text{Err}_2(h)$$

$$\frac{\text{Err}_2(h)}{\|h\|^2} \xrightarrow{h \rightarrow 0} 0$$

$$f: \mathbb{R} \rightarrow \mathbb{R}, \quad f(a+h) \simeq f(a) + f'(a)h + \frac{1}{2} h^2 f''(a)$$

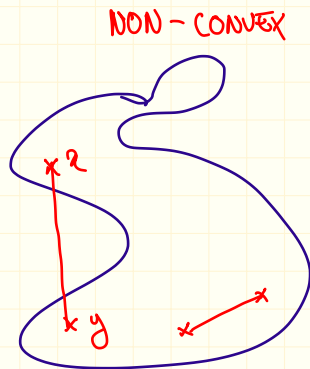
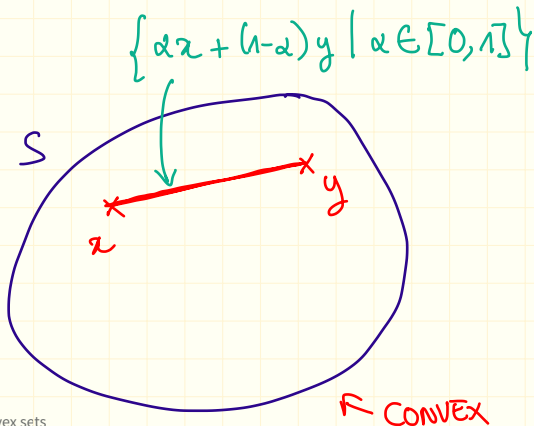
Convex sets

Convex set

Definition

A set $S \subset \mathbb{R}^n$ is called a convex set if for all $x, y \in S$ and all $\alpha \in [0, 1]$,

$$\underline{\alpha x + (1 - \alpha)y \in S}$$

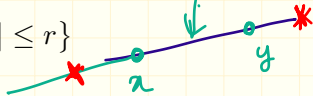


Exercise

1. Show that any subspace S of \mathbb{R}^n is convex.
2. Let $\|\cdot\|$ be a (arbitrary) norm and $r \geq 0$. Show that the "ball" of radius r :

$$B(r) = \{x \in \mathbb{R}^n \mid \|x\| \leq r\}$$

is convex.



① Let $x, y \in S$ and $\alpha \in [0, 1]$, since S is closed under linear combinations, $\alpha x + (1-\alpha)y \in S$. S is convex

② Let $x, y \in B(r)$, $\alpha \in [0, 1]$

$$\|\alpha x + (1-\alpha)y\| \leq \|\alpha x\| + \|(1-\alpha)y\| = \alpha\|x\| + (1-\alpha)\|y\|$$

triangular ineq.

$$\leq \alpha r + (1-\alpha)r = r$$

$$\rightarrow \alpha x + (1-\alpha)y \in B(r).$$

Convex functions

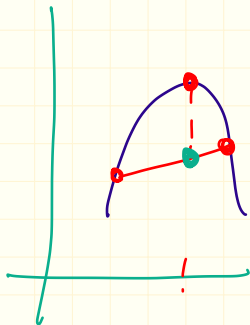
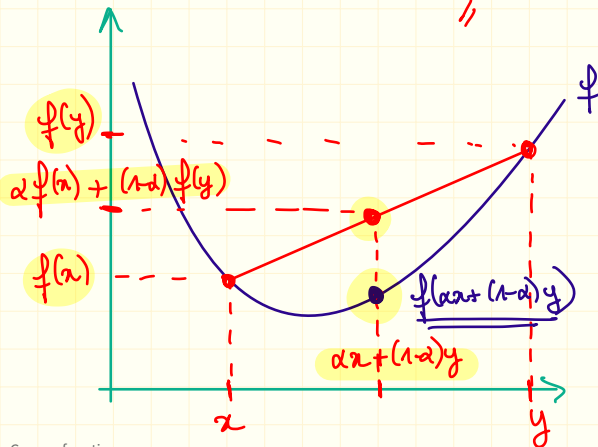
Convex / concave functions

Definition

A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if for all $x, y \in \mathbb{R}^n$ and all $\alpha \in [0, 1]$,

$$\underline{f(\alpha x + (1 - \alpha)y)} \leq \underline{\alpha f(x) + (1 - \alpha)f(y)}. \quad (1)$$

\gg



Convex / concave functions

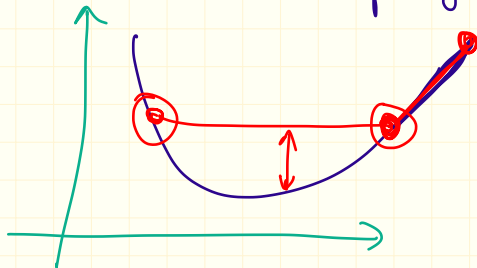
Definition

A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if for all $x, y \in \mathbb{R}^n$ and all $\alpha \in [0, 1]$,

$$f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y). \quad (1)$$

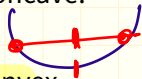
- We say that f is *strictly convex* if there is strict inequality in (1) whenever $x \neq y$ and $\alpha \in (0, 1)$.
- A function f is called *concave* if $-f$ is convex.

→ inequality in the other direction.



Exercise

1. Show that any linear map $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex and concave.
2. Show that a norm $\| \cdot \|$ is convex.
3. Show that the sum of two convex functions is also a convex function.



$$f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y)$$

① let $x, y \in \mathbb{R}^n$ and $\alpha \in [0, 1]$

$$f(\alpha x + (1-\alpha)y) \stackrel{\text{①}}{\leq} \alpha f(x) + (1-\alpha)f(y) \quad \leftarrow \text{because } f \text{ is linear}$$

① \rightarrow f is convex

② \rightarrow f is concave

$$\textcircled{2} \quad \|\alpha x + (1-\alpha)y\| \leq \|\alpha x\| + \|(1-\alpha)y\| = \underbrace{\alpha}_{\geq 0} \|x\| + \underbrace{(1-\alpha)}_{\geq 0} \|y\|$$

\uparrow by trig. inequality.

Convex functions and derivatives

Convex functions vs their tangents

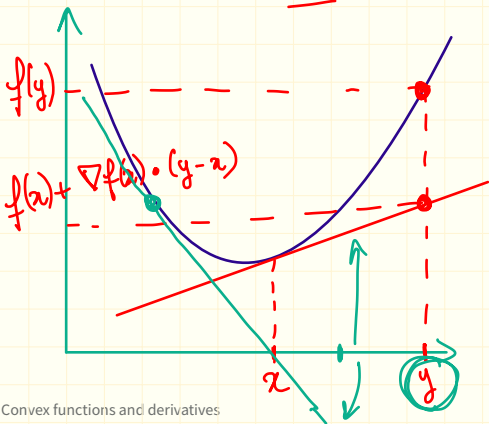
Proposition

A differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if and only if for all $x, y \in \mathbb{R}^n$

$$f(y) \geq f(x) + \langle \nabla f(x), (y - x) \rangle.$$

$\xrightarrow{x+h}$
 $x+h$

\xrightarrow{h}
 h



"The graph of a convex function is above its tangents"

Proof

⊆ Let's assume that for all $x, y \in \mathbb{R}^n$, $f(y) \geq f(x) + \nabla f(x) \cdot (y - x)$

Let $x, y \in \mathbb{R}^n$ and $\alpha \in [0, 1]$, define $z_\alpha = \alpha x + (1 - \alpha)y$.

Applying the inequality twice:

$$\begin{aligned} \alpha f(x) &\geq \alpha f(z_\alpha) + \nabla f(z_\alpha) \cdot (\alpha x - z_\alpha) \\ (1 - \alpha) f(y) &\geq (1 - \alpha) f(z_\alpha) + \nabla f(z_\alpha) \cdot ((1 - \alpha)y - z_\alpha) \end{aligned}$$

SUM: $\alpha f(x) + (1 - \alpha) f(y)$

$$\begin{aligned} &\geq f(z_\alpha) + \nabla f(z_\alpha) \cdot (\alpha x + (1 - \alpha)y - z_\alpha) \\ &= 0 \end{aligned}$$

f is convex

Proof

\Rightarrow let assume that f is convex.

let $x, y \in \mathbb{R}^n$, let $t \in [0, 1]$. write $z_t = (1-t)x + ty$
 $= x + t(y-x)$

Taylor's formula: $f(z_t) = f(x + \overset{h}{t(y-x)})$
 $= \underline{f(x)} + t \nabla f(x) \cdot (y-x) + \text{Err}_1(t)$

f is convex $f(z_t) \leq \underline{(1-t)f(x)} + t f(y)$

Combining: $t f(y) \geq t f(x) + t \nabla f(x) \cdot (y-x) + \text{Err}_1(t)$

$f(y) \geq f(x) + \nabla f(x) \cdot (y-x) + \frac{\text{Err}_1(t)}{t}$

$t \rightarrow 0 \rightarrow 0$

Proof

Minimizers of a convex function

Corollary

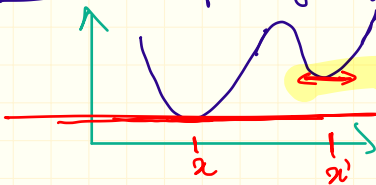
$$f(x+h) \approx f(x) + \nabla f(x) \cdot h + \frac{1}{2} h^T \nabla^2 f(x) h$$

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a differentiable convex function and $x \in \mathbb{R}^n$.

Then

$$\underline{x \text{ is a minimizer of } f} \iff \underline{\nabla f(x) = 0}.$$

Proof: \Rightarrow True for any differentiable f



$$\Leftarrow \text{Let } y \in \mathbb{R}^n, \quad f(y) \geq f(x) + \underbrace{\nabla f(x)}_{=0} \cdot (y-x)$$

$\geq f(x) \rightarrow x \text{ is a minimizer of } f.$

Hessian of convex function

Proposition

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a twice-differentiable function. Then f is convex if and only if for all $x \in \mathbb{R}^n$, $H_f(x)$ is positive semi-definite.

Recall : • a matrix M is PSD if for all $v \in \mathbb{R}^n$,
 $v^T M v \geq 0$

• this is equivalent to saying that all the eigenvalues of M are ≥ 0

Example : $f(x) = \|x\|^2 = x_1^2 + \dots + x_n^2$

For all $x \in \mathbb{R}^n$, $H_f(x) = 2I_n \leftarrow$ PSD hence f is convex

Hessian of convex function

Proposition

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a twice-differentiable function. Then f is convex if and only if for all $x \in \mathbb{R}^n$, $H_f(x)$ is positive semi-definite.

Remarks: if $f : \mathbb{R} \rightarrow \mathbb{R}$ is twice differentiable:

f is convex $\Leftrightarrow f''(x) \geq 0$ for all $x \in \mathbb{R}$

Remark: If $H_f(x)$ is positive definite for all $x \in \mathbb{R}^n$ then f is strictly convex. But the converse is not true in general.

Proof idea: $f(y) \geq f(x) + \nabla f(x) \cdot (y-x) + \frac{(y-x)^T H_f(x) (y-x)}{2} \geq 0$

Jensen's inequality

Jensen's inequality

Theorem

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function. Then for all $x_1, \dots, x_k \in \mathbb{R}^n$ and all $\alpha_1, \dots, \alpha_k \geq 0$ such that $\sum_{i=1}^k \alpha_i = 1$ we have

$f(\text{"average"}) \rightarrow f\left(\sum_{i=1}^k \alpha_i x_i\right) \leq \sum_{i=1}^k \alpha_i f(x_i)$ ← "average of the f 's"
↗ if f is concave

More generally, if X is a random variable that takes value in \mathbb{R}^n we have

$$f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)].$$

Example • $f: \mathbb{R} \rightarrow \mathbb{R}$ is convex therefore,
 $x \mapsto x^2$

for any random variable X , $\mathbb{E}[X]^2 \leq \mathbb{E}[X^2]$

$$\text{Var}(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2 \geq 0$$

Example: entropy

- let's consider a random variable X taking values in $1, \dots, k$
- let's write $p_i = P(X=i)$ for $i \in \underline{1, \dots, k}$.

The entropy of X is defined as:

$$H(X) = \sum_{i=1}^k p_i \log\left(\frac{1}{p_i}\right) \geq 0$$

\log is a concave function (exercise!) hence, by Jensen:

$$H(X) = \sum_{i=1}^k p_i \log\left(\frac{1}{p_i}\right) \leq \log\left(\sum_{i=1}^k p_i \cdot \frac{1}{p_i}\right) = \log(k)$$

Example: entropy

$$\rightarrow 0 \leq H(X) \leq \log(k)$$

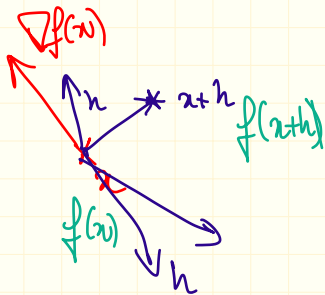
The value $\log(k)$ is achieved for $p_i = \frac{1}{k}$
for all i :

$$H(X) = \sum_{i=1}^k \frac{1}{k} \log\left(\frac{1}{\frac{1}{k}}\right) = \log(k)$$

The entropy is maximal when X is uniformly distributed over $\{1, \dots, k\}$

Questions?

$$f(x+h) \simeq f(x) + \langle h, \nabla f(x) \rangle$$



$$\max \{ \underline{\sigma \cdot x} \mid \|x\| = 1 \}$$

Questions?

$$\# \text{paths}(k+1, i \rightarrow j) = \sum_{\ell \text{ neighbor of } j} \# \text{paths}(k, i \rightarrow \ell)$$

$$\frac{\nabla f(x)}{\|\nabla f(x)\|} = \underset{\substack{\sigma \in \mathbb{R}^n \\ \|\sigma\|=1}}{\operatorname{argmax}}$$

$$\lim_{t \rightarrow 0}$$

$$\frac{f(x+t\sigma) - f(x)}{t}$$