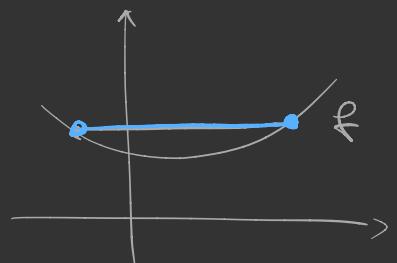


Convex function: connecting segment is always larger than the function itself, i.e.

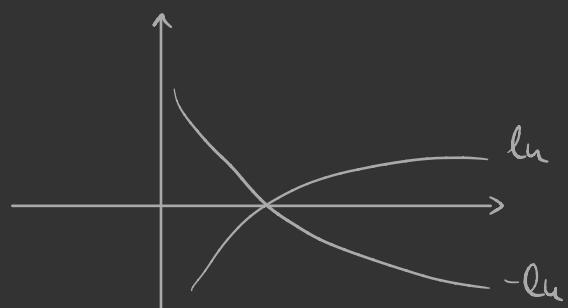


For  $t \in [0,1]$ :

$$f(tx + (1-t)y) \leq t \cdot f(x) + (1-t) \cdot f(y).$$

Examples:  $t \mapsto t^2$ ,  $t \mapsto \exp(t)$

$$t \mapsto -\ln(t)$$



For  $X$  a random variable. Then by Jensen's inequality we can obtain

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

### 3.2 Cramér's Theorem and Hoeffding's inequality.

Connecting the law of large numbers and large deviation bounds.

Recall the moment generating function (Laplace transform):

$$M_X(\theta) = \mathbb{E}[\exp(\theta X)].$$

We define the cumulant generating function

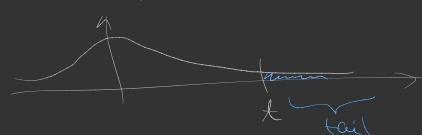
$$C_X(\theta) = \ln(\mathbb{E}[\exp(\theta X)]).$$

### Theorem 3.6: Cramér's Theorem

Let  $X_1, \dots, X_M$  be a sequence of independent random variables with cumulant generating functions  $C_{X_l}$ ,  $l \in [M]$ . Then for  $t > 0$

$$\underbrace{\mathbb{P}\left(\sum_{l=1}^M X_l \geq t\right)}_{\text{probability to "produce a number in the tail" }} \leq \exp\left(\inf_{\theta > 0} \left\{-\theta t + \sum_{l=1}^M C_{X_l}(\theta)\right\}\right).$$

probability to "produce a number in the tail"



Proof: Fix  $\theta > 0$ . Then Markov inequality yields

$$P\left(\sum_{l=1}^M X_l \geq t\right) = P\left(\exp\left(\theta \sum_{l=1}^M X_l\right) \geq \exp(\theta t)\right)$$

Markov  $\leq \frac{\mathbb{E}\left(\exp\left(\theta \cdot \sum_{l=1}^M X_l\right)\right)}{\exp(\theta t)}$

$$\begin{aligned} \exp(a+b) &= \exp(a) \cdot \exp(b) \\ &\Rightarrow \frac{\mathbb{E}\left(\prod_{l=1}^M \exp(\theta X_l)\right)}{\exp(\theta t)} \\ \text{independent } &\Rightarrow \frac{\prod_{l=1}^M \mathbb{E}[\exp(\theta X_l)]}{\exp(\theta t)} \\ C_{X_l}(\theta) &= \ln(\mathbb{E}[e^{\theta X_l}]) \\ \Rightarrow \exp(C_{X_l}(\theta)) &= \mathbb{E}[e^{\theta X_l}] \\ &\Rightarrow \frac{\prod_{l=1}^M \exp(C_{X_l}(\theta))}{\exp(\theta t)} \\ &= \exp\left(-\theta t + \sum_{l=1}^M C_{X_l}(\theta)\right). \end{aligned}$$

Since this is true for every  $\theta > 0$ , it is in particular for the infimum over  $\theta > 0$ .  $\blacksquare$

### Theorem 3.8: Hoeffding's inequality

Let  $X_1, \dots, X_M$  be a sequence of indep. random variables, such that for  $l \in [M]$ :  $\underbrace{\mathbb{E} X_l = 0}$  and  $\underbrace{|X_l| \leq B_l}$  almost surely.

not really a restriction, because we can define  $\tilde{X}_l = X_l - \mathbb{E}[X_l]$ . Means:  $P(|X_l| \leq B_l) = 1$ , does not hold for Gaussian RV.

Then:

$$P\left(\sum_{l=1}^M X_l \geq t\right) \leq \exp\left(-\frac{t^2}{2 \sum_{l=1}^M B_l^2}\right).$$

easy function with known quantities

Concept of "almost surely": emphasise that there might be realizations, which contradict the statement but these occasions are so isolated that they "don't count".

As a consequence (using  $X_e$ ):

$$P\left(\sum_{e=1}^M |X_e| \geq t\right) \leq 2 \exp\left(-\frac{t^2}{2 \sum_{e=1}^M B_e^2}\right) \xrightarrow{t \rightarrow \infty} 0$$

Remark: For  $t > 0$ :  $P(|X| \geq t) = P(X > y) + P(X < -y)$

$\nwarrow$  distinct events  $\nearrow$

Remark: Using basic transformations, we get that

$$(\text{in particular using } t = S \cdot M \text{ and } \left| \frac{1}{M} \sum_{e=1}^M (X_e - E[X_e]) \right| \stackrel{\text{since } M}{\leq} \frac{1}{M} \sum_{e=1}^M |X_e - E[X_e]|)$$

$$P\left(\left|\frac{1}{M} \sum_{e=1}^M (X_e - E[X_e])\right| \leq S\right) \geq 1 - 2 \exp\left(-\frac{M^2 S^2}{2 \sum_{e=1}^M B_e^2}\right)$$

complement

If the  $X_e$ 's are identically distributed, i.e.  $X_e \sim X$ , then  
 $E[X_e] = E[X]$  and  $B_e = B$ . Then we can simplify

$$P\left(\left|\left(\frac{1}{M} \sum_{e=1}^M X_e\right) - E[X]\right| \leq S\right) \geq 1 - 2 \exp\left(-\frac{MS^2}{2B^2}\right).$$

gives our first understanding  
 has many samples  $M$  we need  
 to estimate the mean value

$\leftarrow$  "overwhelming probability"  
 $M \rightarrow \infty \rightarrow 0$ .

which gives us a Law of Large Numbers (LLN).

Intuition about the tail of RV's:

Tails are telling us, how far away a realization is from its expected value.

If our RV is bounded almost surely by  $B$ ,  
 we have a density which is compactly supported by  $[-B, B]$ .

$$\rightarrow E[X] \in [-B, B]$$

$$\rightarrow X - E[X] \leq 2B \text{ a.s.}$$

For RV's for which the tails decay very fast, we need less samples to get a good estimate for the mean value.

# Vorlesung 16

16.06.2020

Proof for claim 3.8:  $P\left(\sum_{l=1}^n X_l \geq t\right) \leq \exp\left(-\frac{t^2}{2 \cdot \sum_{l=1}^n B_l^2}\right)$ .

By assumption,  $X_l \in [-B_l, B_l]$ .



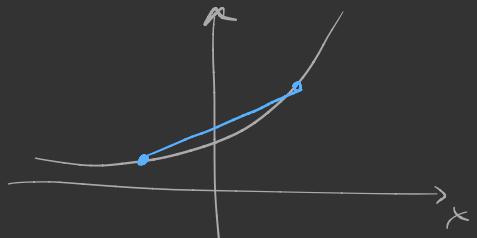
We can write

$$X_l = t(-B_l) + (1-t)B_l \quad \text{for } t \in [0,1].$$

with  $t = \frac{B_l - X_l}{2B_l}$ .

Consider

$$f(x) = \exp(\theta x), \quad \theta > 0.$$



Then  $f$  is convex, hence, by definition, we have for  $t \in [0,1]$

$$\begin{aligned} f(x) &= f(tx_1 + (1-t)x_2) \\ &\leq t \cdot f(x_1) + (1-t) \cdot f(x_2). \end{aligned}$$

Combining this, we get that

$$\begin{aligned} f(X_l) &= \exp(\theta X_l) = f(t(-B_l) + (1-t)B_l) \\ &\leq t \cdot f(-B_l) + (1-t) \cdot f(B_l) \\ &= \frac{B_l - X_l}{2B_l} \exp(-\theta B_l) + \frac{B_l + X_l}{2B_l} \exp(\theta B_l). \end{aligned}$$

Using  $E[X_l] = 0$ , we have

$$E[\exp(\theta X_l)] \leq E\left[\frac{B_l - X_l}{2B_l} \exp(-\theta B_l)\right]$$

$$\leq E\left[\frac{B_l + X_l}{2B_l} \exp(\theta B_l)\right]$$

with  $E\left[\frac{B_l + X_l}{2B_l}\right] \stackrel{\text{lim}}{=} \frac{B_l + E[X_l]}{2B_l} = \frac{B_l}{2B_l} = \frac{1}{2}$ .

$$E[X_l] = 0$$

$$\Rightarrow \mathbb{E}[\exp(\theta X_e)] \leq \frac{1}{2} e^{-\theta B_e} + \frac{1}{2} e^{\theta B_e}$$

Using the series def. of  $e^x$ , we have

$$e^x = \sum_{k=0}^{\infty} \frac{x^k}{k!}.$$

$$\Rightarrow \mathbb{E}[\exp(\theta X_e)] \leq \frac{1}{2} \sum_{k=0}^{\infty} \frac{(-\theta B_e)^k}{k!} + \frac{1}{2} \sum_{k=0}^{\infty} \frac{(\theta B_e)^k}{k!}$$

$$\begin{aligned} & \text{cancel terms with odd exponent} \\ & \text{2x for even exponent} \end{aligned} \Rightarrow \sum_{h=0}^{\infty} \frac{(\theta B_e)^{2h}}{(2h)!}$$

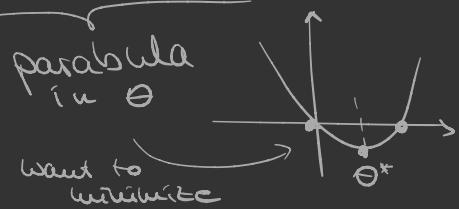
$$(2h)! \geq 2^h \cdot h! \quad \Rightarrow \quad \sum_{h=0}^{\infty} \frac{(\theta B_e)^{2h}}{2^h \cdot h!} = \sum_{h=0}^{\infty} \left( \frac{(\theta B_e)^2}{2} \right)^h \cdot \frac{1}{h!}$$

$$\text{series def.} \quad \Rightarrow \quad \exp\left(\frac{(\theta B_e)^2}{2}\right). \quad \text{by mon. function}$$

$$C_{X_e}(\theta) = \ln \mathbb{E}[\exp(\theta X_e)] \stackrel{\text{Cramer}}{\leq} \ln \left( \exp\left(\frac{(\theta B_e)^2}{2}\right) \right) = \frac{(\theta B_e)^2}{2}$$

$$P\left(\sum_{e=1}^n X_e \geq t\right) \leq \exp\left(\inf_{\theta > 0} \left\{-\theta t + \sum_{e=1}^n C_{X_e}(\theta)\right\}\right)$$

$$\begin{aligned} & \text{exp increasing function} \\ & \leq \exp\left(\inf_{\theta > 0} \left\{-\theta t + \frac{\theta^2}{2} \cdot \sum_{e=1}^n B_e^2\right\}\right) \end{aligned}$$



We know by basic calculations, that

$$\theta^* = \frac{t}{\sum_{e=1}^n B_e^2}$$

$$\Rightarrow P\left(\sum_{e=1}^n X_e \geq t\right) \leq \exp\left(-\frac{t^2}{2 \sum_{e=1}^n B_e^2}\right).$$





Theorem 3.9.

$$\Rightarrow \mathbb{P}(\sum_{e=1}^n X_e > t) \leq 2 \exp\left(-\frac{(kt)^2}{2(2\theta + \theta t)}\right).$$

■

Proof: (Theorem 3.9.) Based on Cramér's theorem.

$$\begin{aligned}\mathbb{E}[\exp(\theta X_e)] &= 1 + \underbrace{\theta \mathbb{E}[X_e]}_{\substack{\uparrow \\ \text{series def.}}} + \sum_{u=2}^{\infty} \frac{\theta^u \mathbb{E}[X_e^u]}{u!} \\ &= 1 + \frac{\theta^2 \theta_e^2}{2} \sum_{u=2}^{\infty} \frac{2\theta^{u-2} \mathbb{E}[X_e^u]}{u! \theta_e^2}\end{aligned}$$

We define  $F_e(\theta) = \sum_{u=2}^{\infty} \frac{2\theta^{u-2} \mathbb{E}[X_e^u]}{u! \theta_e^2}$ . Then we can

$$\text{write } \mathbb{E}[\exp(\theta X_e)] = 1 + \underbrace{\frac{\theta^2 \theta_e^2}{2} \cdot F_e(\theta)}_{= \xi},$$

$$1 + \xi \leq e^\xi \xrightarrow{\quad} \leq \exp\left(\frac{\theta^2 \theta_e^2}{2} \cdot F_e(\theta)\right).$$

Defining  $F(\theta) = \max_{e \in [n]} F_e(\theta)$  and  $\theta^2 = \sum_{e=1}^n \theta_e^2$  we

obtain by Cramér's theorem

$$\ln(\exp(\frac{\theta^2 \theta_e^2}{2} F_e(\theta)))$$

$$\mathbb{P}\left(\sum_{e=1}^n X_e > t\right) \leq \inf_{\theta > 0} \exp\left(-\theta t + \frac{\theta^2 \theta^2 F(\theta)}{2}\right)$$

$$\leq \inf_{0 < R < t} \exp\left(-\theta t + \frac{\theta^2 \theta^2 F(\theta)}{2}\right)$$

Since  $\mathbb{E}[X_e^u] \leq \mathbb{E}|X_e|^u$ , we get

$$F_e(\theta) = \sum_{u=2}^{\infty} \frac{2\theta^{u-2} \mathbb{E}[X_e^u]}{u! \theta_e^2} \leq \sum_{u=2}^{\infty} (R\theta)^{u-2}.$$

↑  
ass.

$$\leq \frac{R\theta^{u-2}}{u! \theta_e^2} \frac{\frac{R^u}{u!} R^{u-2} \theta_e^2}{\cancel{R^{u-2} \theta_e^2}}$$

We know the geometric sum:

$$\sum_{m=0}^{\infty} a^m = \frac{1}{1-a} \quad \text{for } 0 < a < 1.$$

Hence:

$$\sum_{n=2}^{\infty} (R\theta)^{n-2} = \frac{1}{1-R\theta} \quad \text{for } 0 < R\theta < 1.$$

Therefore

$$F(\theta) = \max_l F_l(\theta) \leq \frac{1}{1-R\theta} \quad \text{if } 0 < R\theta < 1$$

and hence

$$P\left(\sum_{k=1}^n X_k > t\right) \stackrel{*}{\leq} \inf_{0 < R\theta < 1} \exp\left(-t\theta + \frac{\theta^2 \theta^2}{2(1-R\theta)}\right)$$

If we choose  $\theta^* = \frac{t}{\theta^2 + Rt}$ , (Ex:  $R\theta^* < 1$ )

then

$$P\left(\sum_{k=1}^n X_k > t\right) \leq \exp\left(-t\theta^* + \frac{(\theta^*)^2 \theta^2}{2(1-R\theta^*)}\right)$$

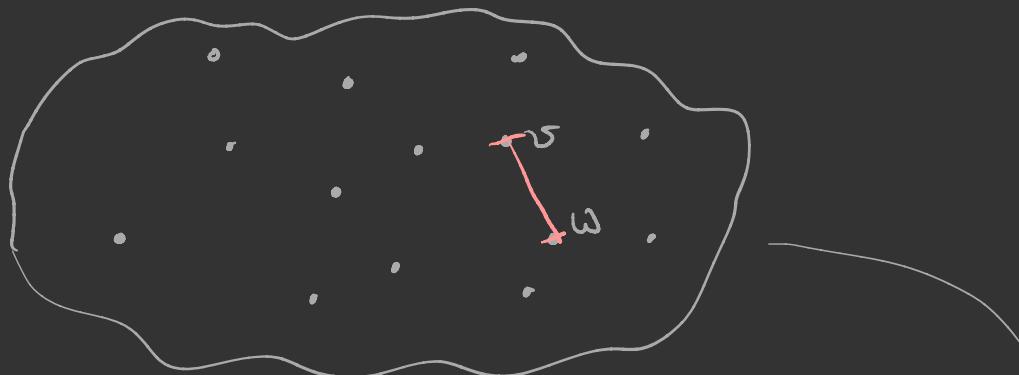
$$\begin{aligned} &= \exp\left(-\frac{t^2}{\theta^2 + Rt} + \frac{t^2}{(\theta^2 + Rt)^2} \cdot \frac{\theta^2}{2(1-R\frac{t}{\theta^2 + Rt})}\right) \\ &\stackrel{\text{basic computations}}{=} \exp\left(-\frac{t^2}{2(\theta^2 + Rt)}\right). \end{aligned}$$

$$= \frac{R\theta^*}{\frac{tR}{\theta^2 + Rt}} < 1$$

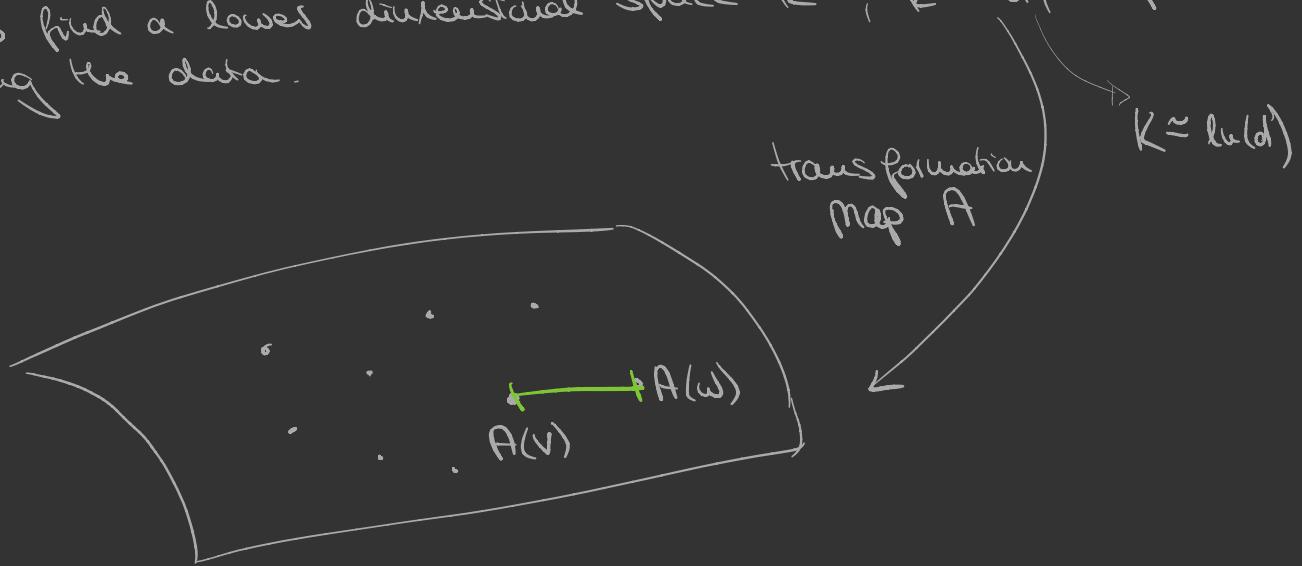
for  $\theta^2 > 0$ .

## 4. Johnson-Lindenstrauss embeddings

Introduction: Point-cloud in high dimension  $\mathcal{C} \subseteq \mathbb{R}^d$ .



Want to find a lower dimensional space  $\mathbb{R}^k$ ,  $k \ll d$ , before processing the data.



Which properties do we want for this embedding / map A?

→ Want to preserve the geometry of the data, meaning keeping the distance between our data points, i.e.

$$\exists \text{ small } \varepsilon > 0 : \forall v, w \in \mathcal{C} \quad (1 - \varepsilon) \|v - w\|_2^2 \leq \|A(v) - A(w)\|_2^2 \leq (1 + \varepsilon) \|v - w\|_2^2$$

↑  
distortion parameter ( $\varepsilon < 1$ )

Since an embedding is called Johnson-Lindenstrauss embedding and will exist. This is a strong example for probabilistic arguments to be beneficial for Data Analysis.



# Vorlesung 17

| 23.06.2020

## Theorem 4.1: Johnson-Lindenstrauss lemma

For any  $0 < \epsilon < 1$  and any integer  $n$ , let  $k \in \mathbb{N}$  such that

$$k \geq 2\beta \left( \frac{\epsilon^2}{2} - \frac{\epsilon^3}{3} \right)^{-1} \quad \begin{matrix} \text{as } \frac{1}{\epsilon^3} \text{ scaling at} \\ k \rightarrow \infty \text{ for } \epsilon \rightarrow 0 \end{matrix}$$

for some  $\beta \geq 2$ . Then for any set  $\mathcal{P}$  of  $n$  points in  $\mathbb{R}^d$ , there is a map  $f: \mathbb{R}^d \rightarrow \mathbb{R}^k$  such that for all  $v, w \in \mathcal{P}$ :

$$(1-\epsilon) \|v-w\|_2^2 \leq \|f(v)-f(w)\|_2^2 \leq (1+\epsilon) \|v-w\|_2^2.$$

Furthermore, the map can be generated at random with probability  $1 - (n^{2-p} - n^{1-p})$ . i.e. not built using knowledge of the data

probability of success  
for the generation

- the larger we choose  $p$ , the higher the success probability but the higher the lower bound for  $k$
- the smaller we choose  $\epsilon$ , the stronger the inequality but the higher the lower bound for  $k$

★ Does  $d$  really not play a significant role here?

## Preliminary:

- ① If  $f$  would be linear map, then the inequality would be equivalent to

$$(1-\epsilon) \leq \|Az\|^2 \leq (1+\epsilon)$$

or

$$|\|Az\|^2 - 1| \leq \epsilon.$$

$$\text{For } z = \frac{v-w}{\|v-w\|_2}.$$

- ② Aim for this map is to estimate the length of a unit vector projected in  $\mathbb{R}^d$  on a randomly picked  $k$ -dim. Subspace.

From a probabilistic view, this is the same as picking the unit vector at random and projecting on a fixed subspace.

$\rightsquigarrow$  we pick an easy subspace  $V = \text{span}\{e_1, \dots, e_k\}$ .

- ③ We generate uniformly distributed vectors  $Y$  on the sphere, using  $X_i \sim N(0,1)$  and  $Y = (X_1, \dots, X_d) / \|X\|_2 \in \mathbb{S}^{d-1}$

- ④ Projection of  $Y$  onto  $V$  is given by  $Z = (y_1, \dots, y_k) \in \mathbb{R}^k$   
i.e. simply taking the first  $k$  coordinates.

$$\Rightarrow \mathbb{E}[\|Z\|_2^2] = \underbrace{\frac{k}{d}}_{\text{Exercise 10.1}} = \mu$$

$$L = \|Z\|_2^2$$

- ⑤ Since we want the length to be close to 1, we can simply rescale  $Z$  accordingly (i.e. later rescale as linear map by this factor).

Summarizing these ideas:

Define  $A = \sqrt{\frac{d}{k}} P_{\text{span}\{e_1, \dots, e_k\}}$

$$X \sim N(0,1)$$

$$Y = \frac{X}{\|X\|_2} \in \mathbb{S}^{d-1}$$

and

$$P_{\text{span}\{e_1, \dots, e_k\}} Y = Z = (y_1, \dots, y_k).$$

$$\text{Exercise: } \mathbb{E}[\|Z\|_2^2] = \frac{k}{d}$$

and hence

$$\mathbb{E}[\|AY\|_2^2] = \mathbb{E}[\|\sqrt{\frac{d}{k}} Z\|_2^2] = 1.$$



By applying Lem. 4.2. a) we get that

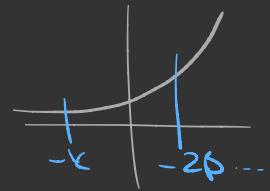
$$\mathbb{P}(L \leq (\underbrace{1-\varepsilon}_{\varepsilon}) \cdot \mu) \leq \exp\left(\frac{\kappa}{2} (1-(1-\varepsilon) + \ln(1-\varepsilon))\right)$$

$$\text{Taylor exp. of } \ln(x), 0 \leq \varepsilon < 1. \quad \ln(1-\varepsilon) \leq -\varepsilon - \frac{\varepsilon^2}{2} \quad \leq \exp\left(\frac{\kappa}{2} (\varepsilon - (\varepsilon + \frac{\varepsilon^2}{2}))\right)$$

$$\begin{aligned} \ln(1-\varepsilon) &\stackrel{\text{Taylor}}{=} \ln(1) + \frac{d}{dx} \ln(x) \Big|_{x=1} \cdot (1-\varepsilon-1) + \frac{1}{2} \frac{d^2}{dx^2} \ln(x) \Big|_{x=1} (1-\varepsilon)^2 + \Theta(\varepsilon^3) \\ &= 0 + \frac{1}{x} \Big|_{x=1} (-\varepsilon) + \frac{1}{2} \cdot (-1) \frac{1}{x^2} \Big|_{x=1} \cdot (-\varepsilon)^2 + \Theta(\varepsilon^3) \\ &= -\varepsilon - \frac{\varepsilon^2}{2} + \Theta(\varepsilon^3). \end{aligned}$$

We have that  $\kappa \geq 2\beta(\frac{\varepsilon^2}{2} - \frac{\varepsilon^3}{3})^{-1} \ln(n)$  by assumption

$$\text{and hence } -\kappa \leq -2\beta(\frac{\varepsilon^2}{2} - \frac{\varepsilon^3}{3})^{-1} \ln(n)$$



$$\Rightarrow \mathbb{P}(L \leq (1-\varepsilon)\mu) \leq \exp\left(\frac{-2\beta \ln(n) \varepsilon^2}{2(\frac{\varepsilon^2}{2} - \frac{\varepsilon^3}{3})}\right)$$

$$= \exp\left(-\frac{\beta \ln(n) \varepsilon^2}{2\varepsilon^2(\frac{1}{2} - \frac{\varepsilon}{3})}\right)$$

$$= \exp\left(-\underbrace{\frac{3}{3-2\varepsilon} \cdot \beta \ln(n)}_{\geq 1} \right) \leq -\beta \ln(n)$$

$$\leq \exp(-\beta \ln(n)).$$

$$= n^{-\beta}.$$

Similarly, we can prove that

$$\mathbb{P}(L \geq (1+\varepsilon)\mu) \stackrel{4.2.(b)}{\leq} \exp\left(\frac{\kappa}{2} (1-(1+\varepsilon) + \ln(1+\varepsilon))\right)$$

$$\begin{aligned} \ln(1+\varepsilon) &\leq \varepsilon - \frac{\varepsilon^2}{2} + \frac{\varepsilon^3}{3} \quad \leq \exp\left(\frac{\kappa}{2} (-\varepsilon + (\varepsilon - \frac{\varepsilon^2}{2} + \frac{\varepsilon^3}{3}))\right) \\ \text{for } \varepsilon > 0 \quad &= \exp\left(-\frac{\kappa(\frac{\varepsilon^2}{2} - \frac{\varepsilon^3}{3})}{2}\right) \end{aligned}$$

$$\begin{aligned} -\kappa &\leq -2\beta(\frac{\varepsilon^2}{2} - \frac{\varepsilon^3}{3})^{-1} \ln(n) \quad \downarrow \\ &\leq \exp(-\beta \ln(n)) = n^{-\beta}. \end{aligned}$$





$$\begin{aligned}
 &= \mathbb{E} \left[ \exp(t(\kappa\alpha - d)(x_1^2 + \dots + x_d^2)) + t\kappa\alpha(x_{k+1}^2 + \dots + x_d^2) \right] \\
 \text{Independence} \downarrow &= \mathbb{E} \left[ \exp(t(\kappa\alpha - d)(x_1^2 + \dots + x_d^2)) \right] \mathbb{E} \left[ \exp(t\kappa\alpha(x_{k+1}^2 + \dots + x_d^2)) \right]
 \end{aligned}$$

$x \sim$  identically

$$= \mathbb{E} \left[ \exp(t(\kappa\alpha - d)x^2) \right]^k \mathbb{E} [t\kappa\alpha x^2]^{d-k}$$

$$= (1 - 2t(\kappa\alpha - d))^{-\frac{k}{2}} (1 - 2t\kappa\alpha)^{-\frac{d-k}{2}} =: g(t).$$

$$\Rightarrow t\kappa\alpha < \frac{1}{2} \quad \text{and} \quad t(\kappa\alpha - d) < \frac{1}{2}$$

$$\Rightarrow 0 < t < \frac{1}{2\kappa\alpha}$$

We want to minimize  $g(t)$  on  $(0, \frac{1}{2\kappa\alpha})$

$$\Leftrightarrow \max_{t \in (0, \frac{1}{2\kappa\alpha})} f(t) = \frac{1}{g(t)^2} = (1 - 2t\kappa\alpha)^{d-k} (1 - 2t(\kappa\alpha - d))^k$$

For this, we try to solve  $f'(t_0) = 0$ . With this we find

$$t_0 = \frac{\frac{1-\alpha}{2\alpha(d-\kappa\alpha)}}{f(t_0)} \in (0, \frac{1}{2\kappa\alpha})$$

$$\left\{ \begin{array}{l} g(t_0) = \sqrt{\frac{1}{f(t_0)}} \quad \text{and} \quad f(t_0) = \left( \frac{d-k}{d-\kappa\alpha} \right)^{d-k} \left( \frac{1}{\alpha} \right)^k. \\ \qquad \qquad \qquad = \alpha^{k/2} \left( 1 - \frac{(1-\alpha)k}{d-k} \right)^{\frac{d-k}{2}} \end{array} \right.$$

$$= \alpha^{k/2} \left( 1 - \frac{(1-\alpha)k}{d-k} \right)^{\frac{d-k}{2}}$$

1

② Exercise:



$$\alpha^{k/2} \left( 1 - \frac{(1-\alpha)k}{d-k} \right)^{\frac{d-k}{2}} \leq \exp \left( \frac{k}{2} (\lambda - \alpha + \ln(\alpha)) \right)$$

(b) can be proven similarly.

|| (a)

Improvement 1: Generate a matrix  $G \in \mathbb{R}^{k \times d}$

To:  $f(x)$  is a Johnson-Lindenstrauss embedding.

Theorem 4.3 : Let  $0 < \varepsilon < \frac{1}{2}$ , any  $n \in \mathbb{N}$ ,  $K \in \mathbb{N}$  st.

$$k \geq \beta \varepsilon^{-2} \ln(n)$$

$$f(\gamma) = \underbrace{\left( \frac{1}{\Gamma(2)} G \right)}_{A} \gamma \quad \text{for} \quad G = (G_{ij})_{i,j}, G_{ij} \sim N(0,1)$$

fulfills  $\forall r, \omega \in \mathcal{B} :$

$$(1-\varepsilon) \|v-w\|_2^2 \leq \|f(v)-f(w)\|_2^2 \leq (1+\varepsilon) \|v-w\|_2^2$$

with probability  $\geq 1 - (n^{2-\beta(1-\varepsilon)} - n^{1-\beta(1-\varepsilon)})$ .

for  $\beta \gg 2$   $\xrightarrow{u \rightarrow \infty} \lambda$

Proof: next time.

# Vorlesung 18

| 24.06.2020

For the proof of Thm. 4.3, we first state the following lemma:

Lemma 4.4: Let  $x \in \mathbb{R}^d$  be fixed. Assume that  $G \in \mathbb{R}^{k \times d}$  has iid entries  $G_{ij} \sim N(0, 1)$ . Then for  $A = \frac{1}{\sqrt{k}} G$ ,

$$\mathbb{P}\left(\left|\|Ax\|_2^2 - \|x\|_2^2\right| > \varepsilon \|x\|_2^2\right) \leq 2 \exp\left(-\frac{(\varepsilon^2 - \varepsilon^3)k}{4}\right)$$

or equivalently

$$\mathbb{P}\left((1-\varepsilon)\|x\|_2^2 \leq \|Ax\|_2^2 \leq (1+\varepsilon)\|x\|_2^2\right) \geq 1 - 2 \exp\left(-\frac{(\varepsilon^2 - \varepsilon^3)k}{4}\right).$$

Proof: (Thm 4.3)

Note again that there are  $\binom{n}{2}$  pairs of points  $v, w \in \mathcal{S}$ . Then

$$\mathbb{P}\left(\text{for some } v, w \in \mathcal{S} : \left|\|f(v) - f(w)\|_2^2 - \|v-w\|_2^2\right| > \varepsilon \|v-w\|_2^2\right)$$

at least one pair  
but can we make  
no intersections  
 $\subseteq \mathbb{P}\left(\bigcup_{v,w} \dots\right)$

$$\leq \sum_{v,w \in \mathcal{S}} \mathbb{P}\left(\left|\|A(v-w)\|_2^2 - \|v-w\|_2^2\right| > \varepsilon \|v-w\|_2^2\right)$$

$$\leq \sum_{v,w \in \mathcal{S}} 2 \exp\left(-\frac{(\varepsilon^2 - \varepsilon^3)k}{4}\right)$$

$k \geq \beta \varepsilon^{-2} \ln(n)$

$$\leq 2 \binom{n}{2} \cdot \exp\left(-\frac{(\varepsilon^2 - \varepsilon^3)\beta \varepsilon^{-2} \ln(n)}{4}\right)$$

$$= 2 \binom{n}{2} n^{-\beta(1-\varepsilon)}$$

$$= n^{2-\beta(1-\varepsilon)} - n^{1-\beta(1-\varepsilon)}.$$

$\Rightarrow$  Complementary probability:  $\forall v, w \in \mathcal{S}$

$$(1-\varepsilon)\|v-w\|_2^2 \leq \|Av - Aw\|_2^2 \leq (1+\varepsilon)\|v-w\|_2^2$$

with probability  $(1 - n^{2-\beta(1-\varepsilon)} - n^{1-\beta(1-\varepsilon)})$ .

Proof: (lemma 4.4.)

$$G_{ij} \stackrel{iid}{\sim} N(0,1), \quad A = \frac{1}{\sqrt{k}} G.$$

1. Step:  $\mathbb{E}\{\|Ax\|_2^2\} = \|x\|_2^2$  for any  $x \in \mathbb{R}^d$  fixed.

let's first check about  $\mathbb{E}\{(Gx)_j^2\}$   $j$ -th coordinate of  $Gx$ .

We recall that  $\mathbb{E}[G_{ij}] = 0$ , moreover  $\mathbb{E}[G_{ij}^2] = 1$ .

$$\begin{aligned} \mathbb{E}\{(Gx)_j^2\} &= \mathbb{E}\left[\left(\sum_{i=1}^d G_{ij} x_i\right)^2\right] \\ &= \mathbb{E}\left[\left(\sum_{i=1}^d G_{ij} x_i\right) \left(\sum_{i=1}^d G_{ij} x_i\right)\right] \\ &= \mathbb{E}\left[\sum_{i=1}^d \sum_{i'=1}^d G_{ii'} G_{ij} x_i x_{i'}\right] \\ &\stackrel{\text{linear}}{=} \sum_{i=1}^d \sum_{i'=1}^d \underbrace{\mathbb{E}\{G_{ii'} G_{ij}\}}_{\text{deterministic}} x_i x_{i'} \\ &\stackrel{\text{indep.}}{=} \sum_{i=1}^d \sum_{i'=1}^d \underbrace{\mathbb{E}\{G_{ii'}\} \mathbb{E}\{G_{ij}\}}_{= \delta_{ii'}} x_i x_{i'} \\ &= \sum_{i=1}^d x_i^2 = \|x\|_2^2. \end{aligned}$$

$$\Rightarrow \mathbb{E}\{\|Ax\|_2^2\} = \frac{1}{k} \sum_{j=1}^k \mathbb{E}\{(Gx)_j^2\} \stackrel{1.\text{ Step}}{=} \frac{1}{k} \sum_{j=1}^k \|x\|_2^2 = \|x\|_2^2.$$

$$\Rightarrow Z_j = \frac{(Gx)_j}{\|x\|_2^2} = \frac{\sum_{i=1}^d G_{ij} x_i}{\|x\|_2^2} \stackrel{1.\text{ Step}}{\sim} N(0, 1).$$

Indeed:

$$\mathbb{E}[Z_j] = \frac{\sum_{i=1}^d \mathbb{E}[G_{ij}] x_i}{\|x\|_2^2} = \underbrace{\sum_{i=1}^d \mathbb{E}[G_{ij}]}_{=0} = 0$$

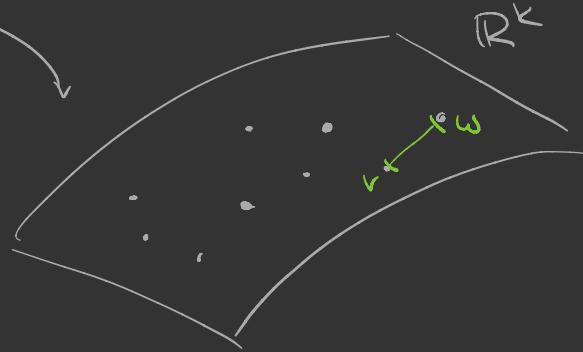
$$\mathbb{E}[Z_j^2] = [1.\text{ Step}] = 1.$$



What we just proved:



$$A = \frac{1}{\sqrt{K}} G$$



with distortion

$$\frac{\|v\|}{\|w\|} \approx (1 \pm \epsilon).$$

Problem: Gaussian R.V. can be unbounded, but very large entries in  $G$  could be problematic for certain applications.

Improvement 2:  $G_{ij} = \begin{cases} 1, & \text{w.p. } \frac{1}{2} \\ -1, & \text{w.p. } \frac{1}{2} \end{cases}$   $\Rightarrow G = \begin{pmatrix} 1 & 1 & 1 \\ 1 & -1 & -1 \\ \dots & \dots & \dots \end{pmatrix}$

and  $A = \frac{1}{\sqrt{K}} G$  will actually give the exact same theorem as Thm 4.3.

Problem:  $G$  might be huge and therefore difficult to store.

Improvement 3:

$$G_{ij} = \begin{cases} \sqrt{\frac{3}{K}}, & \text{w.p. } \frac{1}{6} \\ 0, & \text{w.p. } \frac{2}{3} \\ -\sqrt{\frac{3}{K}}, & \text{w.p. } \frac{1}{6} \end{cases}$$

gives a very sparse matrix, but we can find that

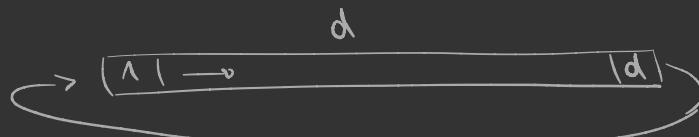
$$\Rightarrow A = \sqrt{\frac{3}{K}} \cdot \begin{pmatrix} -1 & 0 & 0 & 1 & 0 & 0 & \dots \\ 1 & 0 & -1 & 0 & 0 & 0 & \dots \\ 0 & 0 & 0 & -1 & \dots & \dots & \dots \end{pmatrix}^T$$

is much easier to store.

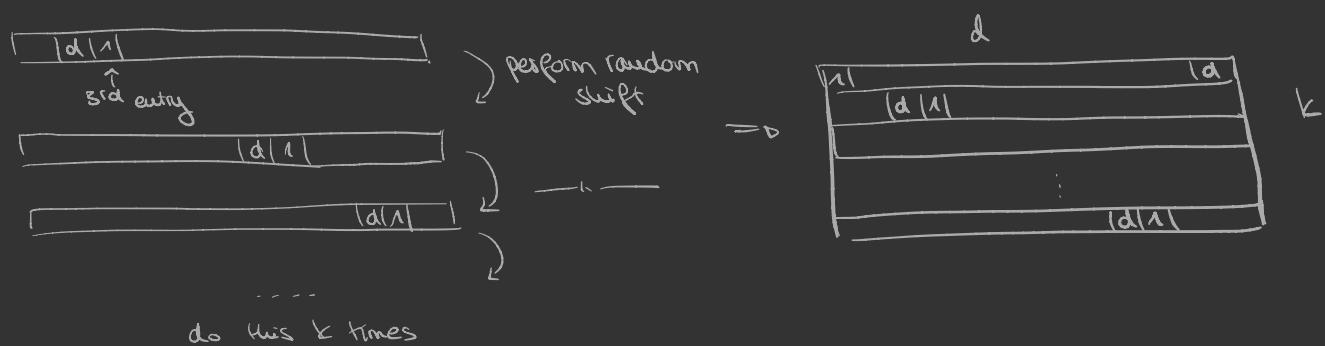
Problem: High randomization, i.e. need to generate and store the full matrix  $\mathbf{G}$  or  $\mathbf{A}$  respectively.

Improvement 4: Use a circulate-random matrix

It starts with vector  $(1, 0, \dots, 0, d) \in \mathbb{R}^d$ , and is constructed using random shifts:



Do this  $k$  times to construct the matrix:



If we use a full circular matrix

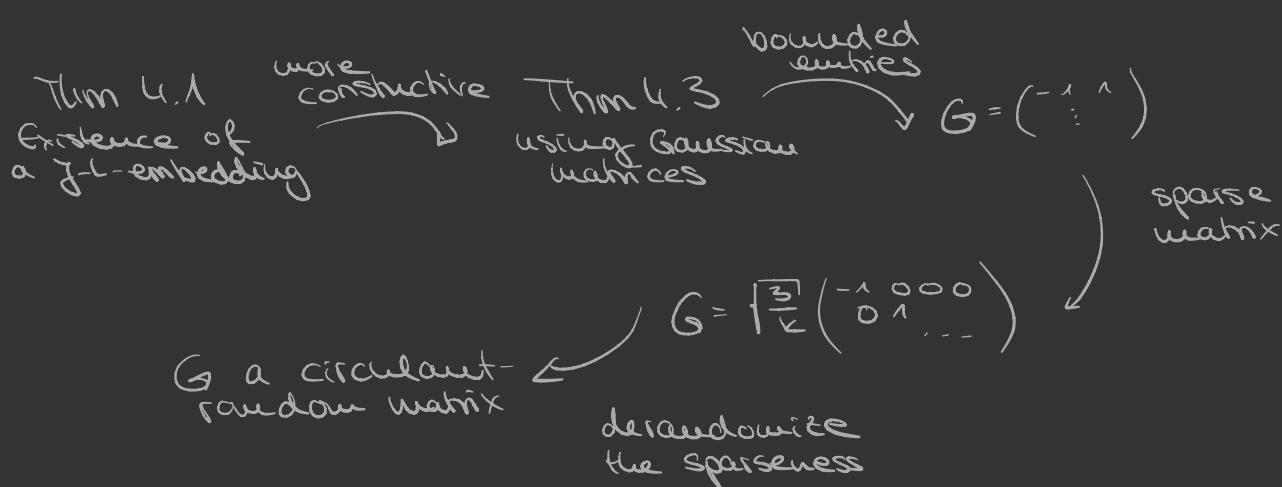
$$\begin{pmatrix} 1 & & & d \\ d & \ddots & & \\ & \ddots & \ddots & \\ & & d & 1 \\ & & & d & 1 \end{pmatrix} \quad = \quad \begin{pmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \end{pmatrix}$$

the matrix vector multiplication can be executed very fast using the Fast-Fourier-Transformation with complexity  $(d \cdot \log(d))$

In a second step we then randomly select  $k$  entries from the result vector (this corresponds with performing  $k$  random shifts).

Such circular-random matrices are actually also a Johnson-Lindenstrauss-embedding.

## Summary of Chapter 4:



But There is no deterministic construction for Johnson-Lindenstrauss embeddings.

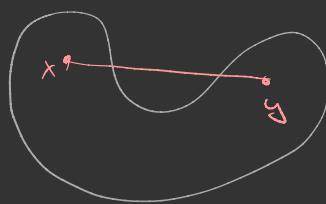
30.06.2020

# Lecture 19

## 5. Convex Analysis

### 5.1 Convex Sets

not convex:



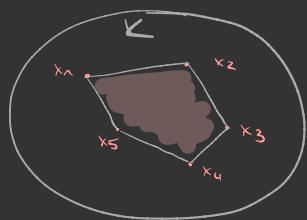
convex:



Def 5.1: A subset  $K \subseteq \mathbb{R}^n$  is called convex if  $\forall x, y \in K$ , the segment between  $x$  and  $y$  is completely contained in  $K$ , i.e.

$$tx + (1-t)y \in K \quad \forall t \in [0,1]$$

Property:  $K \subseteq \mathbb{R}^n$  is convex

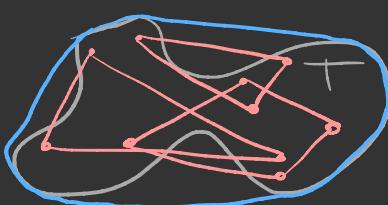


$\Leftrightarrow \forall x_1, \dots, x_n \in K$  and  $t_1, \dots, t_n \geq 0$  s.t.  $\sum_{j=1}^n t_j = 1$ ,  
the convex combination  $\sum_{j=1}^n t_j x_j \in K$ .

Visualization:  $x_1, \dots, x_n$  define vertices and all the space between them is contained in  $K$ .

Def 5.2: Let  $T \subseteq \mathbb{R}^n$  be a set. We define the convex hull of  $T$

as  $\text{Conv}(T) = \left\{ \sum_j t_j x_j : t_j \geq 0, \sum_j t_j = 1, x_j \in T \right\}$ .

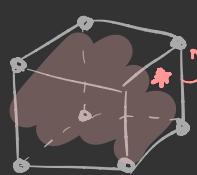


$\text{Conv}(T)$

union of all the polytopes  
we can form using points in  $T$

By definition  $\text{Conv}(T)$  is the smallest convex set, which contains  $T$ .

Property: Often, we need to optimize over non-convex sets,  
i.e. in integer optimizations over  $\{-1,1\}^n \cap \mathbb{Z}^n$ .

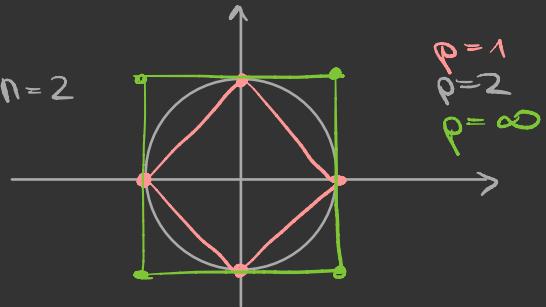


In these situations we can consider a convex relaxation of the problem, in this case optimizing over the domain  $\{-1,1\}^n$  and then find the closest point to the optimal solution of the relaxed problem.

Example:  $B_p(0,1)$ , the ball around 0 with radius 1, w.r.t. the  $p$ -norm.

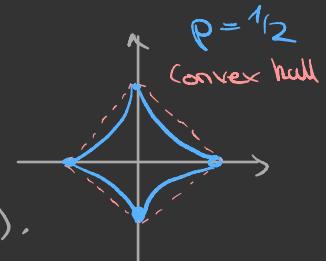
$$\|x\|_p = \left( \sum_{j=1}^n |x_j|^p \right)^{1/p}$$

$$B_p(0,1) = \{ \|x\|_p \leq 1 \}.$$



$B_p(0,1)$  is convex for  $1 \leq p \leq \infty$ .

$B_p(0,1)$  is not convex for  $0 < p < 1$ .



Exercise: For  $0 < p < 1$ :  $\text{Conv}(B_p(0,1)) = B_p(0,1)$ .

More examples:

- Linear subspaces

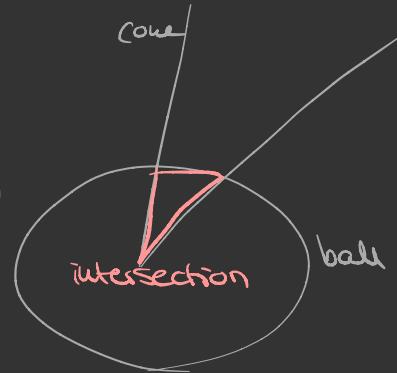
- convex cones

- balls of norms

- intersection of convex sets

- Systems of linear (in)equalities

$$a_1x_1 + \dots + a_nx_n \leq b.$$



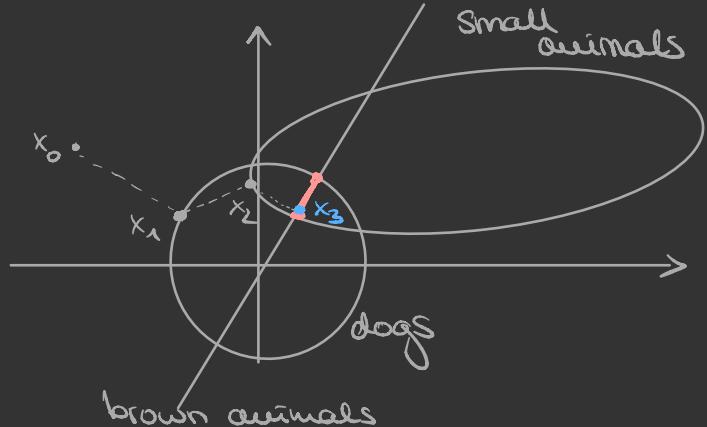
Example: POCS: "Projection onto convex sets"

want to find  $x$ , s.t.

$$x \in C_1 \text{ dog.}$$

$$x \in C_2 \text{ small animals}$$

$$x \in C_3 \text{ brown animals}$$



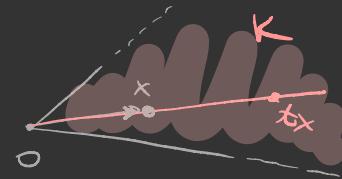
POCS defines an iteration to find such a point:

$$x_{n+1} = P_{C_3} P_{C_2} P_{C_1} x_n.$$

see Exercise 3.2:  
Alternating Projection  
Theorem

Def. 5.3: A set  $K \subseteq \mathbb{R}^n$  is called cone, if

$$\forall x \in K \quad \forall t \geq 0 : \quad tx \in K.$$



If in addition,  $K$  is convex, then  $K$  is called convex cone.

Remark: By definition, since  $t=0$  is allowed,  $O \in K$ .

Exercise: If  $K$  is a convex cone, then

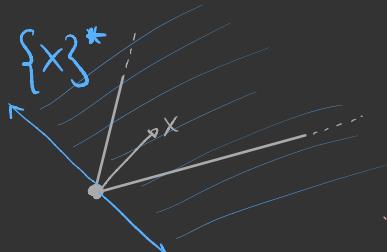


$$\forall x, y \in K \quad \forall s, t \geq 0 : \quad sx + ty \in K.$$

Def.: Let  $K \subseteq \mathbb{R}^n$  be a cone. Then, its dual cone  $K^*$  is defined

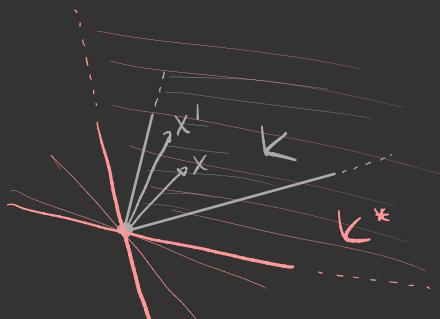
by

$$K^* = \{z \in \mathbb{R}^n \mid \langle x, z \rangle \geq 0, x \in K\}.$$



$\langle x, z \rangle \geq 0$  describes a half space (convex) for a fixed  $x \in K$ .

$$x_1 z_1 + x_2 z_2 + \dots + x_n z_n \geq 0$$



$$\Downarrow \quad \forall x \in K$$

intersection of convex sets.

Exercise: If  $K$  is a closed cone, then  $(K^*)^* = K$ .



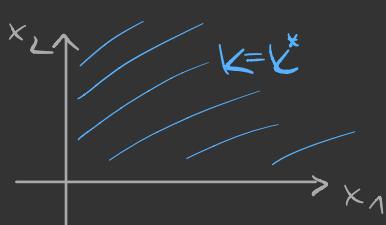
Exercise:  $H, K \subseteq \mathbb{R}^n$  cones. Then

$$H \subset K \Rightarrow K^* \subset H^*.$$

→ "switches" due to more constraints by larger set.

Property: There exist cones which are self-dual, i.e.  $K^* = K$ .  
One example is

$$K = \mathbb{R}_+^n = \{x \in \mathbb{R}^n \mid x_i \geq 0, i \in [n]\}.$$



Theorem 5.4 : Hahn-Banach useful to understand SVMs

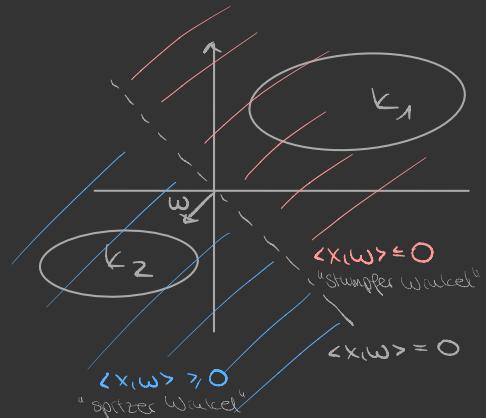
let  $K_1, K_2 \subseteq \mathbb{R}^n$  convex sets, s.t. their interior have an empty intersection, i.e.  $K_1 \cap K_2 = \emptyset$ . Then  $\exists w \in \mathbb{R}^n$  and  $\lambda \in \mathbb{R}$ , s.t.

$$K_1 \subset \{x \in \mathbb{R}^n \mid \langle x, w \rangle \leq \lambda\}$$

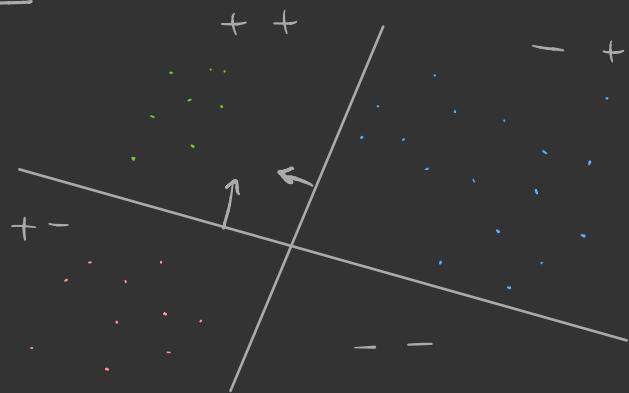
and

$$K_2 \subset \{x \in \mathbb{R}^n \mid \langle x, w \rangle \geq \lambda\}.$$

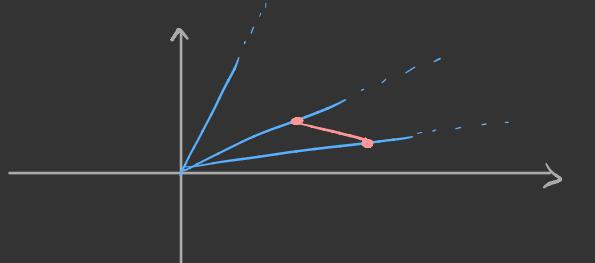
i.e.  $\langle x, w \rangle = \lambda$  separates  $K_1, K_2$ .



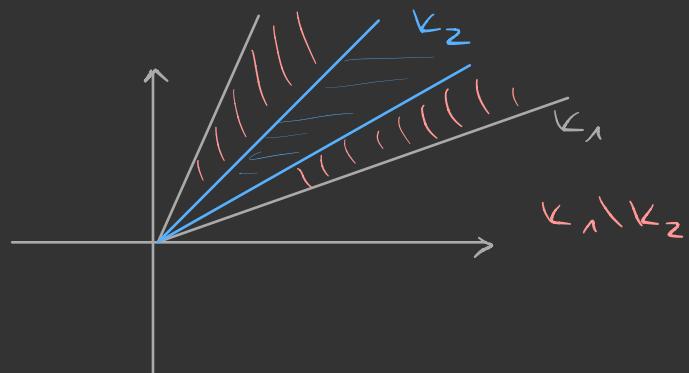
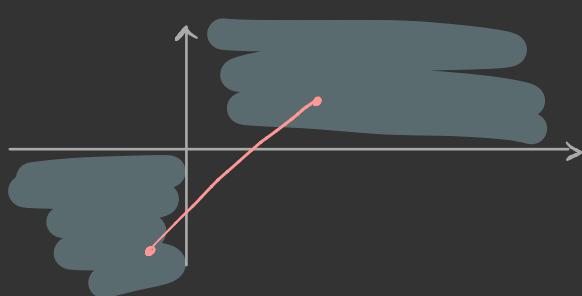
Approach in SVM :



Example for non-convex cone:



- union of lines is a cone
- not convex



## 5.2. Convex Functions

Definition: let  $F: \mathbb{R}^n \rightarrow (-\infty, \infty]$ . We define the domain of  $F$  by

$$\text{dom}(F) = \{x \in \mathbb{R}^n \mid F(x) \neq +\infty\}.$$

A function whose domain is not empty, i.e.  $\text{dom}(F) \neq \emptyset$ , is called proper.  $\rightsquigarrow$  function might be equal to infinity in some points, but not everywhere / there are some function values  $< \infty$ .

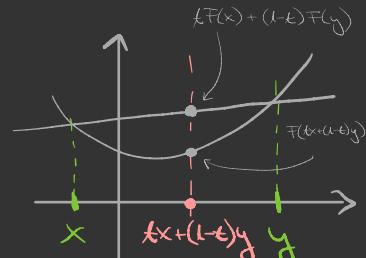
Remark: Let  $F: K \rightarrow \mathbb{R}$ ,  $\emptyset \neq K \subseteq \mathbb{R}^n$ . We can always extend  $F$  to be an extended valued function  $\tilde{F}: \mathbb{R}^n \rightarrow [-\infty, \infty]$  by setting  $\tilde{F}(x) = \infty$  for all  $x \in \mathbb{R}^n \setminus K$ . Then  $\text{dom}(\tilde{F}) = K$ .

Def. 5.7: let  $F: \mathbb{R}^n \rightarrow (-\infty, \infty]$ . Then  $F$  is called convex, if  $\forall x, y \in \mathbb{R}^n$ ,  $t \in [0,1]$ , we have that

$$F(tx + (1-t)y) \leq tF(x) + (1-t)F(y).$$

i.e. the connecting segment is always "over" the graph.

$F$  is called strictly convex, if the inequality is strict for  $t \in (0,1)$ .

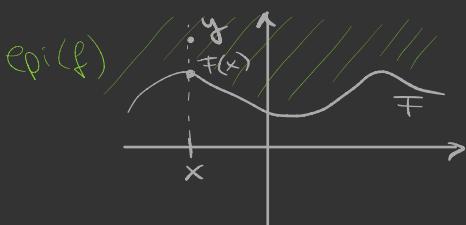


Remark: Define the epigraph of  $F$  by

$$\text{epi}(F) = \{(x,y) \in \mathbb{R}^n \times \mathbb{R} \mid F(x) \leq y\}.$$

Then  $F$  is convex (as a function)

$\Leftrightarrow \text{epi}(F)$  is convex (as a set).



$\rightsquigarrow$  Relationship between convex sets & convex functions.

# Vorlesung 20

01.07.2020

Def:  $F: \mathbb{R}^N \rightarrow (-\infty, \infty]$  is called Strongly convex with parameter  $\gamma > 0$  if for all  $x, y \in \mathbb{R}^N$  and  $t \in [0, 1]$ :

$$F(tx + (1-t)y) \leq tF(x) + (1-t)F(y) - \frac{\gamma}{2} \underbrace{t(1-t)}_{\text{parabola}} \|x-y\|_2^2$$

$F$  is called (strictly, strongly) concave if  $-F$  is (strictly, strongly) convex.

$\sim$  concavity  $\sim$  maximization  
 convexity  $\sim$  minimization

Exercise: Strongly convex  $\Rightarrow$  Strictly convex  $\Rightarrow$  Convex

Exercise: The domain of a convex function is convex (as a set).

Def: A function  $F: K \rightarrow \mathbb{R}$  defined on a convex set  $K \subseteq \mathbb{R}^n$  is called convex if its canonical extension  $\tilde{F}: \mathbb{R}^N \rightarrow (-\infty, \infty]$  is convex.

Proposition 5.8: Let  $F: \mathbb{R}^N \rightarrow \mathbb{R}$  be diffable.

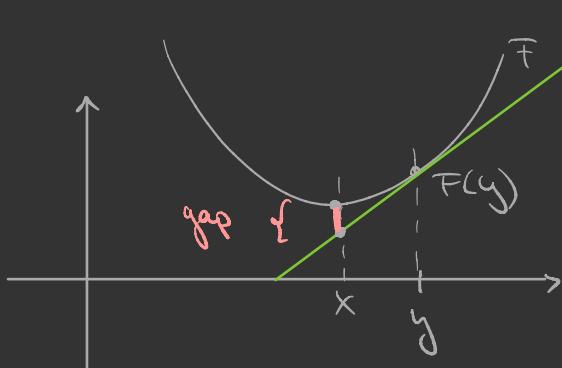
(i)  $F$  is convex

$$\Leftrightarrow F(x) \geq F(y) + \underbrace{\langle \nabla F(y), x-y \rangle}_{\text{linear approximation}}$$

$\sim$  In general convex functions are not smooth



$\forall x, y \in \mathbb{R}^n$



(ii)  $\mathcal{F}$  is strongly convex with parameter  $\gamma > 0$

$$\Leftrightarrow \mathcal{F}(x) \geq \mathcal{F}(y) + \langle \nabla \mathcal{F}(y), x-y \rangle + \frac{\gamma}{2} \|x-y\|_2^2 \quad \forall x, y \in \mathbb{R}^N$$

(iii) If  $\mathcal{F}$  is twice differentiable. Then  $\mathcal{F}$  is convex

$$\Leftrightarrow \nabla^2 \mathcal{F}(x) \succeq 0 \quad \forall x \in \mathbb{R}^n$$

↪ Hessian is pos. Semidefinite

Remark: Gradient of  $\mathcal{F}: \mathbb{R}^N \rightarrow \mathbb{R}$

$$\nabla \mathcal{F}(x) = \begin{pmatrix} \frac{\partial}{\partial x_1} \mathcal{F}(x) \\ \vdots \\ \frac{\partial}{\partial x_n} \mathcal{F}(x) \end{pmatrix}$$

Hessian matrix

$$\nabla^2 \mathcal{F}(x) = \left( \frac{\partial^2}{\partial x_i \partial x_j} \mathcal{F}(x) \right)_{i, j \in [N]}$$

$$= \begin{pmatrix} \frac{\partial^2}{\partial^2 x_1} \mathcal{F}(x) & \frac{\partial^2}{\partial x_1 \partial x_2} \mathcal{F}(x) & \cdots & \frac{\partial^2}{\partial x_1 \partial x_n} \mathcal{F}(x) \\ \frac{\partial^2}{\partial x_2 \partial x_1} \mathcal{F}(x) & \frac{\partial^2}{\partial^2 x_2} \mathcal{F}(x) & \cdots & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \frac{\partial^2}{\partial^2 x_n} \mathcal{F}(x) \end{pmatrix}$$

A matrix  $M \in \mathbb{R}^{N \times N}$  is called positive semi-definite,  
short  $M \succeq 0$

$$\Leftrightarrow \forall v \in \mathbb{R} : v^T M v = \langle v, Mv \rangle \geq 0$$

Proof : Taylor expansion  
(idea)

and hence exactly the definition of being convex.

Prop. 5.9:

## Examples

Proposition 5.11: Let  $F: \mathbb{R}^N \rightarrow \mathbb{R}$  (in particular  $\text{dom}(F) = \mathbb{R}^N$ ) be convex. Then  $F$  is continuous.

Remark: No function which takes on the value  $\infty$  can be continuous.  $\rightarrow$  for these we need the following concept

Def. 5.12: A function  $F: \mathbb{R}^N \rightarrow (-\infty, \infty]$  is called lower semicontinuous if for all  $x \in \mathbb{R}^N$  and every sequence  $(x_j)_{j \in \mathbb{N}} \subset \mathbb{R}^N$  with  $x_j \rightarrow x$ , it holds that

$$\liminf_{j \rightarrow \infty} F(x_j) \geq F(x).$$

Example:  $\circ)$   $\chi_K$  if  $K$  is closed, convex.

Remark: A function  $F$  is lower-semicontinuous  $\Leftrightarrow \text{epi}(F)$  is closed.

Def.: A (global) minimizer of a function  $F: \mathbb{R}^N \rightarrow (-\infty, \infty]$  is a point  $x \in \mathbb{R}^N$  s.t.  $\forall y \in \mathbb{R}^N : F(x) \leq F(y)$ .

A local minimum of  $F$  is a point  $x \in \mathbb{R}^N$  s.t.  $\exists \varepsilon > 0$  s.t.  $F(x) \leq F(y) \quad \forall y \in B_\varepsilon(x)$ .

Prop 5.13: Let  $F: \mathbb{R}^N \rightarrow (-\infty, \infty]$  be convex.

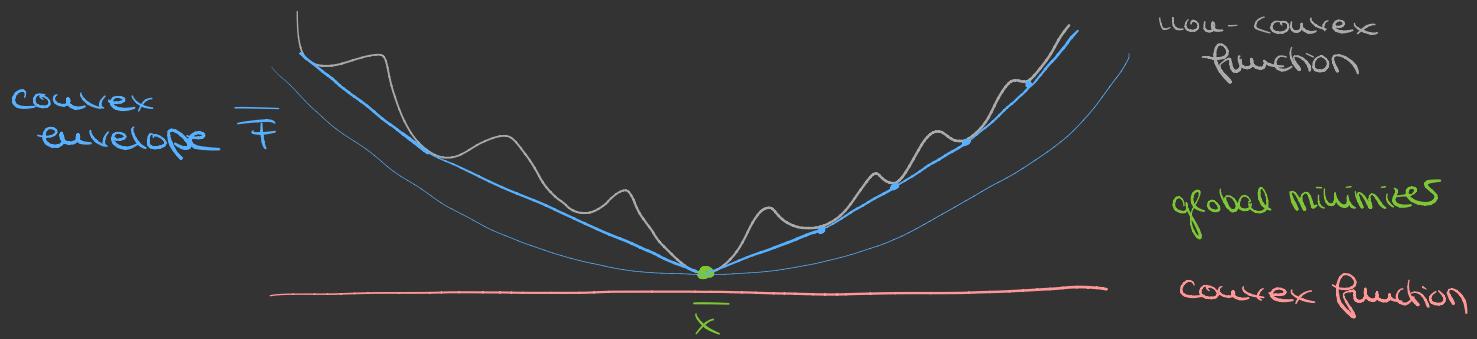
- (i) A local minimum of  $F$  is a global minimum.
- (ii) The set of minima of  $F$  is convex.  $\rightsquigarrow$  equal function value
- (iii) If  $F$  is strictly convex, then the minimum is unique.

Proof (of iii): Assume that  $F$  is strictly convex and  $\exists x, \tilde{x} \in \mathbb{R}^N$ ,  $x \neq \tilde{x}$ , global minimizers, i.e.  $F(x) = F(\tilde{x}) \leq F(y) \quad \forall y \in \mathbb{R}^N$ .

Then

$$\begin{aligned} F(tx + (1-t)\tilde{x}) &\leq tF(x) + (1-t)F(\tilde{x}) = F(x)(t+1-t) \\ &= F(x) \\ &= \min_z F(z) \quad \begin{array}{l} \swarrow \\ \text{global minimizers} \end{array} \\ \Rightarrow x &= \tilde{x}. \end{aligned}$$

## Revisiting the idea of convex relaxation



$$\bar{F}(x) = \max_{G \leq F} G(x) \quad \rightsquigarrow \quad \bar{x} = \underset{x}{\operatorname{argmin}} \bar{F}(x) = \underset{x}{\operatorname{argmin}} F(x)$$

$\bar{x}$

Depending on  $F$ , getting to  $\bar{F}$  can be easy. Also a "good enough" estimation can be sufficient. Sometimes it's as hard as solving the original problem in the first place.

## Another option: Graduating Non-Convexity

We define

$$F(x) = \bar{F}(y) + \lambda \|x - y\|_2^2$$

might want to move adding a parabola which is a strong convex component  
in a "non-convex way"

Then

$$y \mapsto \bar{F}(y) + \lambda \|x - y\|_2^2$$

is convex. With this we can write

$$x^{u+1} = \underset{y}{\operatorname{argmin}} \underbrace{\bar{F}(y) + \lambda_n \|x^u - y\|^2}_{\text{convex}}$$

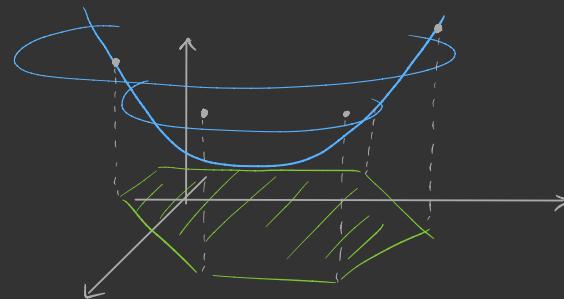
As soon as  $\|x^{u+1} - x^u\| \rightarrow 0$ , this means that we are approaching a minimizer, because  $\lambda_n \|x^u - x^{u+1}\|$  is getting more and more irrelevant in our optimizing steps.

$$\Leftrightarrow \nabla \bar{F}(x^{u+1}) + 2\lambda_n (x^u - x^{u+1}) = 0$$

$$\Leftrightarrow 2\lambda_n (x^{u+1} - x^u) = \nabla \bar{F}(x^{u+1})$$

$$\Leftrightarrow x^{u+1} = x^u - \frac{1}{2\lambda_n} \nabla \bar{F}(x^{u+1}) \quad \text{"gradient descent"}$$

Theorem 5.15: Let  $K \subset \mathbb{R}^N$  be a compact, convex set, and  $F: K \rightarrow \mathbb{R}$  be a convex function. Then  $\bar{F}$  attains its maximum at an extreme point of  $K$ .



Example: Finding 1-rank-matrices in subspaces which are generated by 1-rank-matrices.

$\{u_1, \dots, u_m\}$  orthonormal vectors. We compute

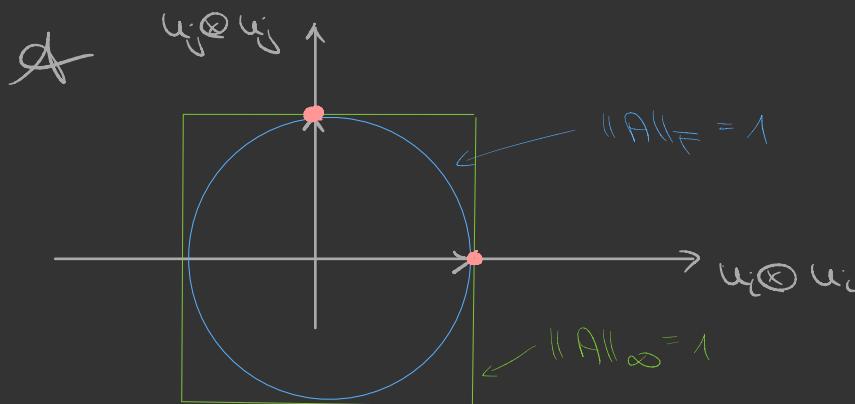
$$u_j \otimes u_j = u_j u_j^T \quad \text{1-rank-matrix.}$$

We define

$$\mathcal{A} = \text{span} \{ u_j \otimes u_j \mid j \in [m] \}.$$

Then  $\forall A \in \mathcal{A}$ :

$$A = \sum_{j=1}^m \alpha_j u_j u_j^T = \sum_{j=1}^m |\alpha_j| \text{sign}(\alpha_j) u_j u_j^T$$



Then  $\langle u_i \otimes u_i, u_j \otimes u_j \rangle_F = \langle u_i, u_j \rangle^2 = \delta_{ij}$  (numerically  $\perp$ )

and hence

1-rank matrix we want to find

Thm 5.15

$$\left\{ u_j \otimes u_j \right\} = \underset{\substack{A \in \mathcal{A} \\ \|A\|_F \leq 1}}{\operatorname{argmax}}$$

convex set

$\|A\|_F$ convex function

any matrix  $A \in \mathcal{A}$ .

$\Rightarrow$  maxima at extremal points / boundary.



# Vorlesung 21:

07.07.2020

## 5.3. Convex Conjugate

Def 5.16: Let  $F: \mathbb{R}^n \rightarrow (-\infty, \infty]$ . Then we define its convex conjugate  $F^*: \mathbb{R}^n \rightarrow (-\infty, \infty]$  by

$$F^*(y) := \sup_{x \in \mathbb{R}^n} \left\{ \underbrace{\langle x, y \rangle - F(x)}_{\text{linear in } y} \right\}.$$

Remark:  $F^*$  is always a convex function.

Property: Fenchel inequality

By definition we have that

$$\langle x, y \rangle \leq F(x) + F^*(y) \quad \forall x, y \in \mathbb{R}^n.$$

Example:  $0 \leq \|x - y\|_2^2 = \|x\|_2^2 + \|y\|_2^2 - 2 \langle x, y \rangle$

$$\Rightarrow \langle x, y \rangle \leq \frac{1}{2} \|x\|_2^2 + \frac{1}{2} \|y\|_2^2$$

Defining  $F(x) = \frac{1}{2} \|x\|_2^2$ , we get that

$$\langle x, y \rangle \leq F(x) + F(y).$$

We can prove that  $F^*(x) = \frac{1}{2} \|x\|_2^2 = F(x)$ .

Exercise:  $1 \leq p < \infty$ . If we consider

$$F(x) = \frac{1}{p} \|x\|_p^p, \quad F^*(x) = ?$$



Proposition 5.17: Let  $F: \mathbb{R}^n \rightarrow (-\infty, \infty]$ . Then

(i)  $F^*$  is lower semicontinuous.

$$F(x_n) \xrightarrow{n \rightarrow \infty} x \quad \liminf_n F^*(x_n) \geq F^*(x).$$

(ii) The biconjugate  $F^{**}$  is the largest, lower semicontinuous function satisfying

$$F^{**} \leq F(x).$$

If  $F$  is convex and lower semicontinuous, then  $F^{**} = F$ .

