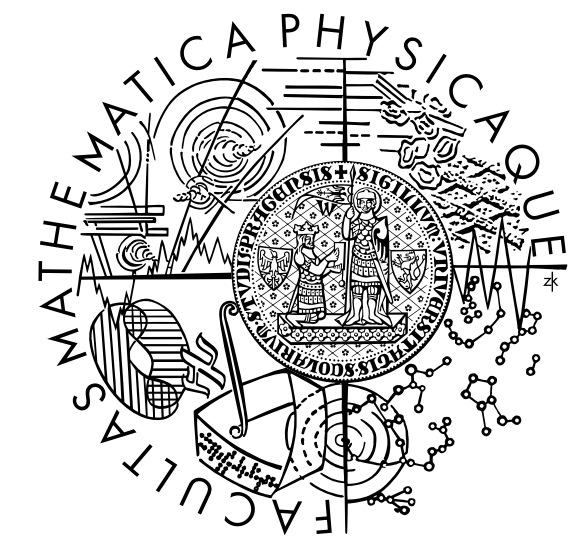


# VÁŽENÁ HLOUBKA DAT A JEJÍ VLASTNOSTI

ONDŘEJ VENCÁLEK

vencalek@karlin.mff.cuni.cz  
Matematicko-fyzikální fakulta Univerzity Karlovy v Praze



**Tento příspěvek z oblasti neparametrických metod analýzy mnohorozměrných dat představuje možnost zobecnění poloprostorové hloubky bodu, tzv. váženou poloprostorovou hloubku, a zabývá se jejími vlastnostmi v závislosti na volbě váhy v její definici. Speciálně nás zajímá nulovost hloubky bodů mimo nosič rozdělení.**

## POLOPROSTOROVÁ HLOUBKA A JEJÍ ZOBECNĚNÍ

Jeden z nejdůležitějších neparametrických přístupů práce s mnohorozměrnými daty je založen na tzv. hloubce dat. Jde vlastně o způsob, jak uspořádat prvky vícerozměrného prostoru. Zřejmě nejrozšířenější definici hloubky publikoval v roce 1975 J. Tukey:

**Definice:** Uvažujme  $p$ -rozměrný prostor  $\mathbb{R}^p$  a bod  $\mathbf{x} \in \mathbb{R}^p$ . Necht  $P$  je pravděpodobnostní míra na  $\mathbb{R}^p$ . Poloprostorová hloubka bodu  $\mathbf{x}$  vzhledem k  $P$  je definována vztahem

$$D_P(\mathbf{x}) = \inf_{\mathbf{u}: \|\mathbf{u}\|=1} P\{\mathbf{y} : \mathbf{u}^T(\mathbf{y} - \mathbf{x}) \geq 0\}$$

Tedy poloprostorová hloubka bodu  $\mathbf{x}$  je definována jako minimální pravděpodobnost uzavřeného poloprostoru obsahujícího tento bod. Navrhované zobecnění této definice, tedy vážená poloprostorová hloubka, místo pravděpodobnosti poloprostoru uvažuje váhovou funkci integrovanou na tomto poloprostoru podle příslušné pravděpodobnostní míry.

**Definice 2:** Uvažujme  $p$ -rozměrný prostor  $\mathbb{R}^p$  a bod  $\mathbf{x} = (x_1, \dots, x_p) \in \mathbb{R}^p$ .

Necht  $P$  je pravděpodobnostní míra na  $\mathbb{R}^p$ .

Bud  $w_+ : \mathbb{R}^p \rightarrow [0, \infty)$  ohraničená, měřitelná váhová funkce, taková, že  $w_+(\mathbf{x}) = 0$  jestliže  $x_p < 0$ .

Necht dále  $w_-(\mathbf{x}) = w_-(x_1, \dots, x_{p-1}, x_p) = w_+(x_1, \dots, x_{p-1}, -x_p)$ .

Váženou poloprostorovou hloubku bodu  $\mathbf{x}$  vzhledem k  $P$  definujeme jako

$$D_{\mathbf{X}}(\mathbf{x}) = D_P(\mathbf{x}) = \inf_{\mathbf{A} \in O_p} \frac{\mathbf{E} P w_+(\mathbf{A}(\mathbf{X} - \mathbf{x}))}{\mathbf{E} P w_-(\mathbf{A}(\mathbf{X} - \mathbf{x}))} \quad (1)$$

kde  $O_p$  označuje prostor všech ortogonálních  $p \times p$  matic.

Definujeme ještě hodnotu podílu dvou nul  $0/0 = 1$ .

Příklady volby váhové funkce pro dvourozměrný prostor:

- $w_+(\mathbf{x}) = 1$  pro  $x_2 \geq 0, x_1 \in \mathbb{R}$ ,
- $w_+(\mathbf{x}) = 1$  pro  $x_2 \geq 0, |x_1| < k$ ,
- $w_+(\mathbf{x}) = 1$  pro  $x_2 \geq 0, |x_1| < kx_2$ ,
- $w_+(\mathbf{x}) = 1$  pro  $x_2 \geq 0, |x_1| < k/x_2$ ,
- $w_+(\mathbf{x}) = \Phi(x_1)$  pro  $x_2 \geq 0, x_1 \in \mathbb{R}$ ,

kde  $k$  je nezáporná konstanta,  $\Phi(\cdot)$  distribuční funkce standardního normálního rozdělení  $N(0, 1)$ . Vždy  $w_+(\mathbf{x}) = 0$  pro všechna ostatní  $\mathbf{x}$ . Viz schéma vyznačující oblasti, kde je váhová funkce nenulová.

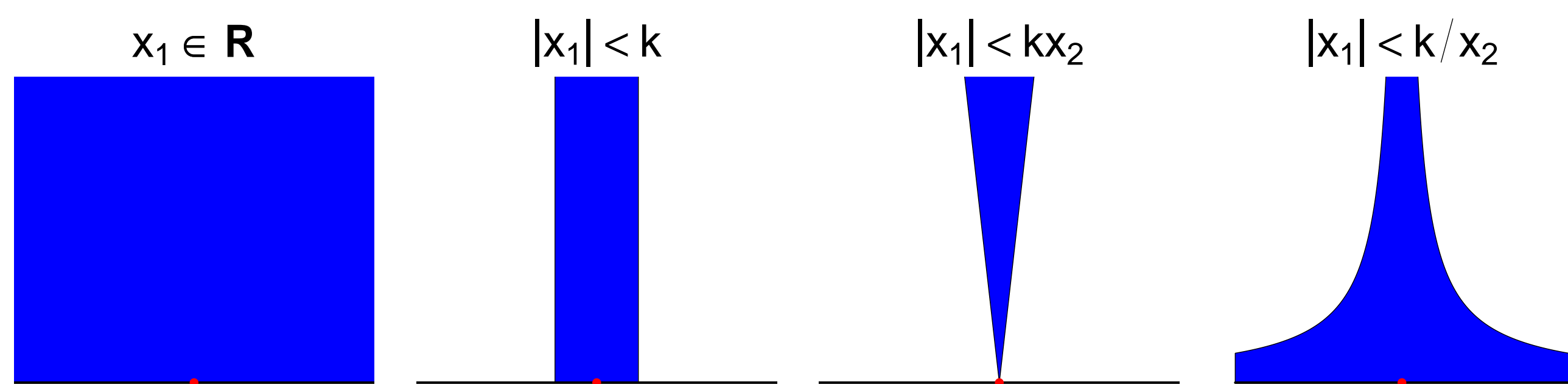


Schéma: příklady oblastí, kde je váhová funkce  $w_+$  nenulová.

## NULOVOST HLOUBKY V BODECH MIMO NOSIČ ROZDĚLENÍ

Jednou z přirozených vlastností, kterou od hloubkové funkce očekáváme, je její nulovost v bodech ležících mimo konvexní uzávěr nosiče rozdělení  $P$ . Připomeňme, že nosič rozdělení  $P$  (značíme jej  $\text{sp}(P)$ ) je nejmenší uzavřená množina s pravděpodobností 1, tedy

$$\text{sp}(P) = \bigcap \{F \in \mathcal{F} : P(F) = 1\},$$

kde  $\mathcal{F}$  je množina všech uzavřených podmnožin  $\mathbb{R}^p$ . Symbolem  $\text{csp}(P)$  pak označujeme uzavřený konvexní obal nosiče rozdělení  $\text{sp}(P)$ . Pro názornost budeme uvažovat dvourozměrný prostor  $\mathbb{R}^2$ , všechny následující věty však mají jednoduché zobecnění pro vícerozměrný případ  $\mathbb{R}^p$ , kde  $p \in \mathbb{N}$ .

Uvažujme nejprve případ, kdy nosič rozdělení je konvexní. Podmínku na váhovou funkci, která je postačující pro to, aby body mimo nosič tohoto rozdělení měly nulovou hloubku, formuluje následující věta:

**Věta 1:** Necht  $\text{sp}(P)$  je konvexní. Označme  $W = \{\mathbf{y} : w_+(\mathbf{y}) > 0\}$ . Předpokládejme, že pro váhovou funkci platí:

$$\forall \mathbf{x} : x_1 = 0, x_2 > 0 \exists U \text{ okolí bodu } \mathbf{x} : U \subset W, \quad (2)$$

potom pro všechny  $\mathbf{x} \notin \text{sp}(P)$  platí  $D_P(\mathbf{x}) = 0$ .

**Důkaz:** V důkazu se využije invariance vážené poloprostorové hloubky vzhledem k posunutí a otočení. Uvažujme bod  $\mathbf{x}_0 \notin \text{sp}(P)$ . Chceme ukázat, že jeho vážená hloubka bude při splnění podmínky (2) rovna nule.

Označme  $\mathbf{x}_m = \arg \min \{\|\mathbf{x} - \mathbf{x}_0\| : \mathbf{x} \in \text{sp}(P)\}$ , tedy bod nosiče rozdělení  $P$  nejbližší k bodu  $\mathbf{x}_0$ . Existence a jednoznačnost tohoto bodu vyplývá z konvexity  $\text{sp}(P)$ . Díky invarianci hloubky vzhledem k posunutí můžeme bůno předpokládat, že  $\mathbf{x}_0 = (0, 0)$ , a díky invarianci hloubky vzhledem k otočení můžeme bůno předpokládat, že  $\mathbf{x}_m = (0, v)$ , kde  $v = \|\mathbf{x} - \mathbf{x}_0\|$ .

Nyní si stačí uvědomit, že pro takovéto otočení musí  $\text{sp}(P) \subset H_v \subset H_0$ , kde  $H_v = \{\mathbf{x} : x_2 \geq v\}$ ,  $H_0 = \{\mathbf{x} : x_2 \geq 0\}$ . Z předpokladu (2) vyplývá, že existuje  $U_{\mathbf{x}_m}$  okolí bodu  $\mathbf{x}_m$  takové, že  $U_{\mathbf{x}_m} \subset W$  a tedy  $U_{\mathbf{x}_m} \cap W \cap \text{sp}(P) \neq \emptyset$ , odtud  $\mathbf{E} P w_+(\mathbf{A}(\mathbf{X} - \mathbf{x}_0)) > 0$ .

Z vlastnosti  $\text{sp}(P) \subset H_v \subset H_0$  naopak vyplývá, že  $\mathbf{E} P w_-(\mathbf{A}(\mathbf{X} - \mathbf{x}_0)) = 0$  a tedy  $D_P(\mathbf{x}_0) = 0$ .  $\square$

Pro nekonvexní nosič rozdělení  $P$  je třeba buď předpokládat souvislost  $\text{sp}(P)$  nebo zesílit požadavek na tvar váhové funkce:

**Věta 2:** Necht existuje  $n \in \mathbb{N}$  takové, že  $\text{sp}(P) = \bigcup_{i=1}^n K_i$ , kde  $K_i$  ( $i=1, \dots, n$ ) je souvislá podmnožina  $\mathbb{R}^2$ , přičemž  $\text{sp}(P)$  nemá žádný izolovaný bod. Bud dále  $m_{ij} = \min \{\|\mathbf{x} - \mathbf{y}\| : \mathbf{x} \in K_i, \mathbf{y} \in K_j\}$ ,  $i, j = 1, \dots, n$ . Uvažujme  $m = \max_{1 \leq i, j \leq n} m_{ij}$ . Necht pro váhovou funkci  $w_+$  platí

$$\forall \mathbf{x} : |x_1| \leq m/2 \exists U_{\mathbf{x}} \text{ okolí } \mathbf{x} \text{ takové, že } U_{\mathbf{x}} \subset W. \quad (3)$$

Potom  $\mathbf{x} \notin \text{csp}(P) \Rightarrow D_P(\mathbf{x}) = 0$ .

Uvědomme si, že pro speciální případ  $n = 1$ , tedy situaci, kdy nosič je souvislý nekonvexní, věta říká, že podmínka (2) zaručuje, že všechny body mimo konvexní obal nosiče rozdělení mají nulovou váženou poloprostorovou hloubku.

## HLOUBKA BODŮ, KTERÉ JSOU V KONVEX. OBALU NOSIČE ROZDĚLENÍ, ALE MIMO NOSIČ SAMOTNÝ

Snadno se ukáže, že poloprostorová hloubka všech bodů v konvexním obalu nosiče rozdělení je nenulová (tedy i v bodech, které se nacházejí v konvexním obalu nosiče rozdělení, ale mimo nosič samotný). Vážená poloprostorová hloubka tuto diskutabilní vlastnost obecně nemá.

Uvažujme například rovnoměrné rozdělení na kruhové výseči vzniklé z kružnice se středem  $S$  a poloměrem  $r$  (viz obrázek 1a). Všechny body v konvexním obalu nosiče rozdělení (viz obrázek 1b) mají nenulovou poloprostorovou hloubku. Uvažujme nyní váženou poloprostorovou hloubku s váhovou funkcí ve tvaru

$$w_+(\mathbf{x}) = 1 \quad \text{pro } |x_1| < h, x_2 \geq 0 \\ = 0 \quad \text{jinak,} \quad (4)$$

kde  $h$  je nějaká kladná konstanta menší než  $r$ . Pak oblast bodů s nenulovou váženou poloprostorovou hloubkou je sjednocením nosiče rozdělení a kružnice se středem  $S$  a poloměrem  $h$  (viz obrázek 1c).

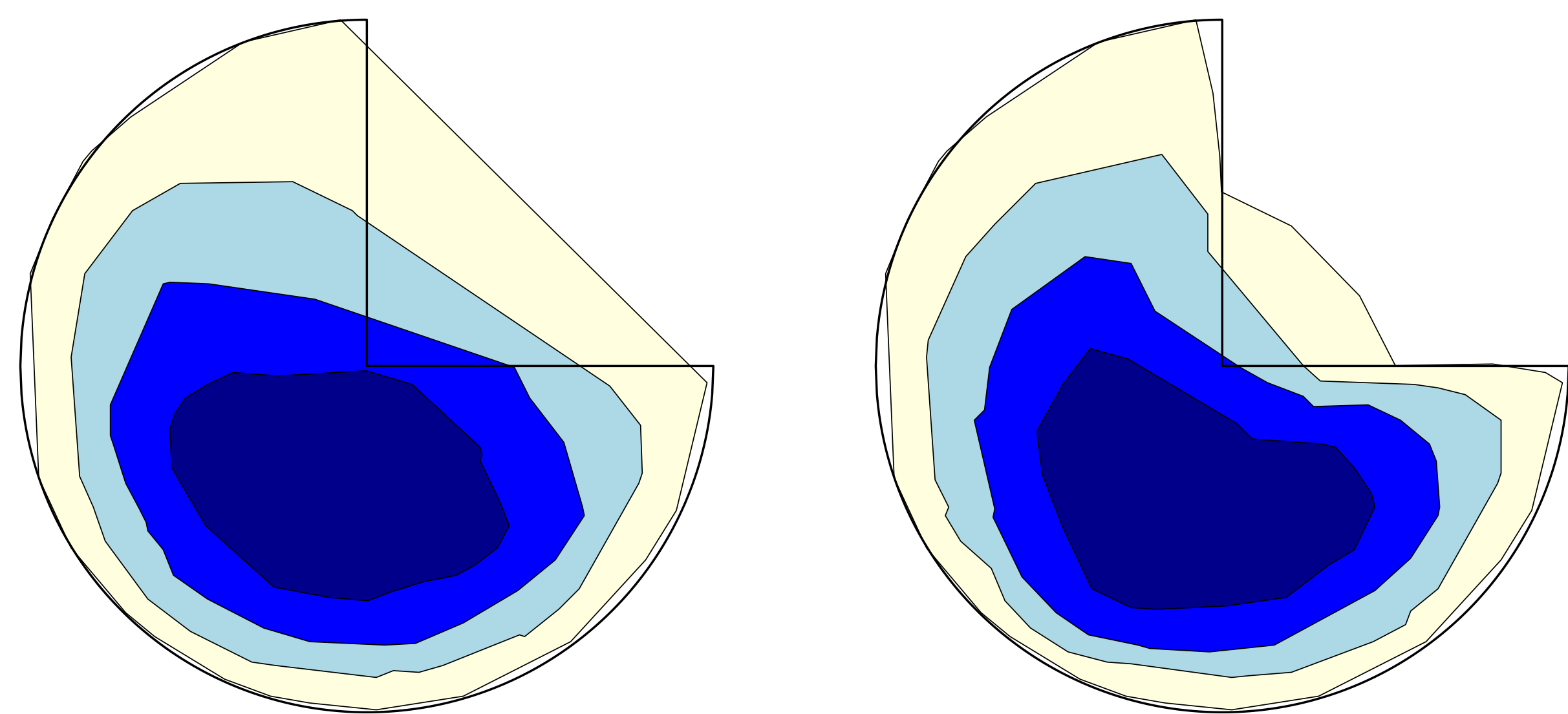


Obrázek 1: nosič rovnoměrného rozdělení ve tvaru kruhové výseče (a); konvexní obal tohoto nosiče je oblastí bodů s nenulovou poloprostorovou hloubkou (b); oblast bodů s nenulovou váženou poloprostorovou hloubkou při použití váhové funkce (4), kde  $h = r/2$  (c).

Výše uvedené skutečnosti ilustrujeme pomocí simulace. Bylo generováno 1000 bodů z rovnoměrného rozdělení na kruhové výseči (poloměr kruhu  $r$ ). Srovnáme-li tvar oblastí 25, 50, 75 a 100% nejhlubších bodů pro poloprostorovou hloubku a váženou poloprostorovou hloubku s váhovou funkcí (4), viz obrázek 2, odhalíme velké rozdílnosti. Jednou z nejdůležitějších je nulovost hloubky velkého množství bodů v oblasti, která se nachází v konvexním obalu nosiče rozdělení, ale mimo nosič samotný, pro váženou poloprostorovou hloubku.

### Poloprostorová hloubka

### Vážená poloprostorová hloubka



Obrázek 2: Empiricky zjištěné oblasti 25, 50, 75 a 100% nejhlubších bodů pro poloprostorovou (vlevo) a váženou poloprostorovou (vpravo) hloubku. Nosič rozdělení, ze kterého bylo generováno, je znázorněn silnější čarou.

**Poděkování:** Rád bych poděkoval docentu Danielu Hlubinkovi a Lukáši Kotíkoví, se kterými na práci o hloubce dat dlouhodobě spolupracuji. Tato práce byla podporována grantem GAČR 201/08/0486.

### Literatura:

- [1] Liu R.Y., Parelius J.M. a Singh K. (1999) *Multivariate analysis by data depth: descriptive statistics, graphics and inference*. The Annals of Statistics, Vol. 27, No. 3, 783–858.
- [2] Zuo Y. a Serfling R. (2000) *General notion of statistical depth function*. Annals of Statistics, Vol. 28, 461–482.
- [3] Tukey J. (1975) *Mathematics and picturing data*. In Proceedings of the 1975 International Congress of Mathematics, Vol. 2, 523–531.