

## Výpočetní teorie strojového učení

PAC učení 1

## Koncept = pojem

- Pojmy často **vysvětlujeme** prostřednictvím příkladů: dobrý politik, krásný člověk, nudná přednáška, ..
- Definice potřebujeme kvůli vzájemné komunikaci*
- Mějme **definiční obor X** popsány jako množinu všech možných instancí objektů, kde objekt je charakterizován hodnotami svých **atributů**.
- Koncept** na množině objektů **X** je definován tak, že každý prvek **X** je označen
  - buď jako příklad daného pojmu,
  - nebo o něm víme, že uvažovanému pojmu **neodpovídá** (nereprezentuje jej).
 Koncept vlastně představuje nějakou podmnožinu **X**, její charakteristická funkce **c** se někdy označuje jako **cílová funkce (target f.)**.  
 Obecně může být **c** libovolná bool.funkce nad **X**, tj. **c: X → {0,1}**

## Hypotéza pro klasifikované příklady

- Naučit se koncept z příkladů** znamená vlastně naučit se booleovskou funkci **c** (na celém definičním oboru **X**) na základě **příkladů** jejího chování na nějaké podmnožině **X**, tedy na základě znalostí o chování na množině **Y ⊂ X**.
- Nechť **D** označuje aktuální množinu **trénovacích příkladů** chování ve tvaru **<x, c(x)>**, což znamená, že **D** je rozdělena na pozitivní a negativní příklady
- Hypotéza** = pokus o popis cílového konceptu.
- H** nechť označuje množinu všech možných *přípustných hypotéz*. Obecně hypotéza **h** může být rovněž booleovská funkce nad **X**, tj. **h: X → {0,1}**. Jinými slovy hypotéza také reprezentuje podmnožinu **X**.

**Cíl učení:** najít hypotézu, která je **korektní** pro všechny příklady z **X**, tj. platí **h(x) = c(x)** pro všechna **x ∈ X**.

## Induktivní bias

**Věta o ošklivém kačátku** pro konečný def.obor. Nechť **D** je klasifikovaná trénovací množina pro koncept **K**, který tvoří podmnožinu *konečného* definičního oboru **X** všech myslitelných možností, které lze vyjádřit v uvažovaném jazyce pro popis trénovacích dat.

Pokud **D ⊂ X** a **D ≠ X**, pak pro každý prvek **y ∈ X \ D** platí, že pravděpodobnost tvrzení „**y** patří ke konceptu **K**“ je rovna **0,5**. Tedy všechny objekty z množiny **D \ E** mají stejnou pravděpodobnost, že do konceptu **K** patří (i že do něj nepatří).

Cílem strojového učení je najít **korektní popis konceptu** na základě omezeného počtu příkladů (množina výrazně menší než **X**). *Lze toho vůbec dosáhnout? Za jakých podmínek a jak postupovat?*

- Využití doménové znalosti pro volbu **tvaru hypotéz**, která zaručí, že algoritmus učení nalezne hypotézu, která je korektní = **induktivní bias** (předpoklad).

## Charakteristika prostoru (množiny) hypotéz

- Příklad:** Pojem „*střední postava*“.  
Zde hypotézou může být libovolný **obdélník rovnoběžný s osami**, tj. podinterval intervalu „*výška X váha*“, tj. **<150, 200> X <50, 120>**.  
Taková hypotéza nemůže popsat libovolnou podmnožinu def.oboru. **Musíme mít tedy dobrý důvod pro volbu takového omezení!**
- Říkáme, že **množina hypotéz H plně pokrývá** (shatters) **množinu dat obsahující M bodů**, pokud pro libovolnou z **2<sup>M</sup>** možných klasifikačních úloh na uvažované množině dat existuje korektní hypotéza **h z H**, která tuto úlohu řeší (dokonale odděluje poz. a neg. příklady v uvažované úloze).
- Kapacitu** dané množiny hypotéz **H** charakterizuje **maximální počet bodů**, které mohou být touto množinou hypotéz **H** pokryty. Tato hodnota se také označuje jako **Vapnik-Chervonenkis (VC) dimenze** množiny hypotéz **H**.
- Např. **obdélník rovnoběžný s osami** v rovině může oddělit 4 body, které nejsou v přímce. VC takové množiny hypotéz je 4. To je velmi málo. Je takový prostor hypotéz vůbec užitečný?

## Co se lze naučit, je-li def.obor nekonečný?

Platí obdoba „Věty o ošklivém kačátku“ i v případě definičního oboru, který **není** konečný?

**Co když uvažujeme nekonečný definiční obor, např. N (přirozená čísla)?**

V takovém případě dokonce musí existovat koncepty, které **nelze popsat pomocí předem zvoleného jazyka pro popis hypotéz!**

**Důvod?** Množina **C** všech konceptů (= množina všech podmnožin **N**) je **výrazně mohutnější** než množina **H** hypotéz (= mn. formulí ve zvoleném jazyce)!

*Presněji:* neexistuje žádné vzájemně jednoznačné zobrazení mezi **C** a **H**

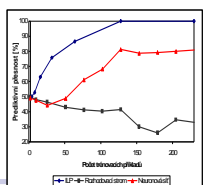
## Důsledky věty o ošklivém kačátku

Upřesnění cíle strojového učení:

Nemá smysl pátrat po *přesně správné hypotéze*, ale hledáme

**skoro správnou hypotézu** (approximately correct), která splňuje *doplňkové požadavky* – Occamova břitva, bias, ..., vhodný kompromis mezi paměťovými nároky pro reprezentaci konceptu a mezi pravděpodobností chybné klasifikace.

Lze nějak charakterizovat chování známé hypotézy pro cílový koncept, jehož popis pro celý definiční obor není k dispozici? „**Křivka učení**“ a testovací data.

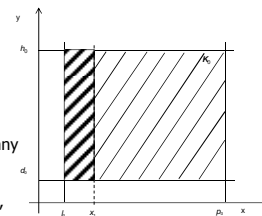


## Algoritmus A,

který se skoro přesně učí pojem „středně velký objekt“:

Pro všechny vstupující pozitivní příklady sledujte hodnoty  $MAX_i$  a  $MIN_i$  ve všech atributech  $i$  pro dosud přečtené vzorky. Hypotéza necht' je pak interval  $\langle MIN_i, MAX_i \rangle$   $X \langle MIN_i, MAX_i \rangle$

**Tvrzení:** K tomu, aby se tento algoritmus naučil hypotézu skoro přesně, stačí mu prohlédnout  $4/\epsilon \cdot \ln(4/\delta)$  příkladů.



Důkaz:

Pro navržený algoritmus **A** platí „všechny generované hypotézy jsou konzistentní, tj. pro libovolnou navrženou hypotézu **h** platí,  $h \subseteq k$ “, kde  $k$  je cíl.koncept.

Tuto vlastnost budeme dále označovat jako **Inkluze**.

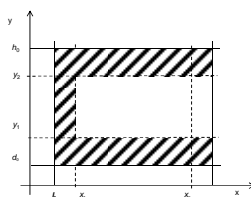
## Důsledek Inkluze:

Pro libovolnou hypotézu **h** navrženou algoritmem **A** platí, že pravděpodobnost jevu „nový objekt **x** bude špatně klasifikován“ je menší než tato pravděpodobnost pro cílový koncept  $k = \langle l, p \rangle \times \langle d, h \rangle$ , ozn.  $P(k)$ .

Jaká je pravděpodobnost toho, že nový objekt **x** bude špatně klasifikován?

Jistě platí, že  $P(x: h(x) \neq k(x)) < P(k)$ .

- Necht'  $P(k) < \epsilon$ , pak jistě i pravděpodobnost chyby je menší než  $\epsilon$ .
- Necht'  $P(k) > \epsilon$ .



Zvolme  $x_0$  tak, že  $x_0 = \inf\{x: P(\langle l, x \rangle \times \langle d, h \rangle) > \epsilon/4\}$

Jistě platí  $P(\langle l, x_0 \rangle \times \langle d, h \rangle) > \epsilon/4$ ,

a tedy pravděpodobnost, že vzorek padne mimo levý sloupec je menší než  $(1-\epsilon)/4$

## Odhad počtu trénovacích příkladů pro koncept vyjádřený intervalem

Pro  $m$  různých vzorků platí:

pravděpodobnost toho, že každý vzorek padne mimo levý sloupec, je **menší** než

$$(1 - \epsilon/4)^m$$

Tentýž postup lze použít i pro ostatní rohy:

Pravděpodobnost jevu „žádný prvek z trénovací množiny obsahující  $m$  vzorků nepadnul do rohu“ je menší než

$$4(1 - \epsilon/4)^m < \delta$$

$$4(1 - \epsilon/4)^m \cong 4e^{-m \cdot \epsilon/4} < \delta$$

$$m > 4/\epsilon \ln(4/\delta)$$

## Východiska výpočetní teorie PAC stroj.učení

Výchozí předpoklad - **stacionarita**: Trénovací i testovací množina jsou vybírány z téže populace za použití totožné distribuce pravděpodobnosti

(Probably Aproximately Correct) **PAC učení**: Kolik trénovacích příkladů je třeba, aby se podařilo eliminovat všechny velmi špatné hypotézy?

- X** množina všech možných příkladů (s distribucí **D**)
- f** skutečný popis konceptu
- H** množina všech možných hypotéz,  $h \in H$  je aktuální hypotéza
- er(h)** = pravděpodobnost jevu „ $x \in X$  a platí  $f(x) \neq h(x)$ “  
=  $P(\{x: x \in X \text{ s distribucí } D \text{ a platí } f(x) \neq h(x)\})$   
tuto hodnotu studují „křivky učení“

Hypotéza **h** je  **$\epsilon$ -skoro správná**, pokud  $er(h) < \epsilon$

## Základní otázka PAC

Bud' **m** mohutnost trénovací množiny

*Můžeme určit **m** tak, aby pouhá konsistence hypotézy s trénovací množinou byla dostatečnou zárukou toho, že jsme našli skoro správnou hypotézu?*

Takový odhad může sloužit např. jako vodítko při shromažďování či posuzování trénovacích dat

V dalším předpokládáme, že hledaný popis konceptu **f** je **prvkem uvažované množiny hypotéz H** (tedy, že existuje korektní, tj.konzistentní a úplná hypotéza pro daný koncept)

### Základní pojmy PAC

**H** konečný prostor všech hypotéz  
**f** (hledaný) cílový koncept z mn. **H**

$H_\epsilon$  je  $\epsilon$  okolí **f**, tj. množina obsahující pouze ty hypotézy, jejichž pravděpodobnost chyby je  $< \epsilon$

$H_{\epsilon\text{-špatná}} = H - H_\epsilon$   
 obsahuje právě ty hypotézy, jejichž pravděpodobnost chyby je větší než  $\epsilon$

**Pokusme se odhadnout, za jakých okolností platí, že**  
 $P(\text{„nalezena hypotéza konzistentní se všemi trén. příklady a současně z } H_{\epsilon\text{-špatná}} \text{“}) < \delta$

### Odhad potřebného počtu příkladů

**Tvrzení:** Necht' **h** je hypotéza konzistentní se všemi trénovacími příklady.  
 Pokud platí  $P(h \in H_{\epsilon\text{-špatná}}) < \delta$ , pak pravděpodobnost toho, že „**h** je  **$\epsilon$ -skoro správná**“ je větší než  $(1 - \delta)$ .

**Zdůvodnění:**  
 Předpokládáme, že v **H** existuje nějaká hypotéza konzistentní se všemi trénovacími příklady. Jistě platí  
 $P(h \in H - H_{\epsilon\text{-špatná}}) + P(h \in H_{\epsilon\text{-špatná}}) = 1$

Víme, že  $H_\epsilon = H - H_{\epsilon\text{-špatná}}$ , a proto  
 $P(h \in H_\epsilon) \geq (1 - \delta)$ .

**h**  $\in H_\epsilon$  znamená, že **h** klasifikuje prvky z **X** s chybou menší než  $\epsilon$ , tj. **h** je  **$\epsilon$ -skoro správná**

### Kdy pro **h** konzistentní s trén. daty platí $P(h \in H_{\epsilon\text{-špatná}}) < \delta$ ?

Necht' **b** je libovolná opravdu špatná hypotéza, tedy h., jejíž pravděpodobnost chyby  $er(b) = \text{pravděpodobnost jevu „ } x \in X \text{ a platí } f(x) \neq h(x) \text{“} > \epsilon$ . V tomto případě je **b**  $\in H_{\epsilon\text{-špatná}}$  a platí:

- Pravděpodobnost opačného jevu „**b** správně klasifikuje jeden náhodně zvolený příklad“ je  $P(b \text{ správně klasifikuje 1 zvolený příklad}) \leq (1 - \epsilon)$
- Pravděpodobnost, že **b** správně klasifikuje **m** náhodně zvol. příkladů:  $P(b \text{ správně klasifikuje } m \text{ zvolených příkladů}) \leq (1 - \epsilon)^m$

Pravděpodobnost, že existuje prvek z  $H_{\epsilon\text{-špatná}} = \{b_1, \dots, b_m\}$ , který správně klasifikuje **m** zvol. příkladů, je rovna pravděpodobnosti, že „**b**, správně klasifikuje **m** zvolených příkladů“ nebo „**b**, správně klasifikuje **m** zvolených příkladů“, tj.  
 $P(h \in H_{\epsilon\text{-špatná}}) \leq \sum_{i=1}^m (1 - \epsilon)^m = |H_{\epsilon\text{-špatná}}| (1 - \epsilon)^m \leq |H| (1 - \epsilon)^m$

Pokud  $|H| (1 - \epsilon)^m < \delta$ , pak jistě  $P(h \in H_{\epsilon\text{-špatná}}) < \delta$

### Jak lze zajistit, aby $P(h \in H_\epsilon) \geq 1 - \delta$ ?

**Stačí, aby platilo  $P(h \in H_{\epsilon\text{-špatná}}) < \delta$**   
**To dosáhneme, když  $|H| (1 - \epsilon)^m < \delta$**

*Fakt:* Pro  $|\epsilon| < 1$ , platí, že  $(1 - \epsilon)^m \approx e^{-\epsilon m}$

Podmínku lze tedy přepsat do tvaru  
 $\ln(|H| e^{-\epsilon m}) < \ln \delta$  čili  $\ln |H| - \epsilon m < \ln \delta$

Postačující podmínka pro počet příkladů **m** je tedy  
 $m \geq (\ln |H| - \ln \delta) / \epsilon = 1/\epsilon * (\ln |H| + \ln (1/\delta))$

Máme-li k dispozici alespoň  $1/\epsilon * (\ln |H| + \ln (1/\delta))$  příkladů a učící algoritmus navrhuje hypotézu **h**, která je se všemi těmito příklady konzistentní, pak pravděpodobnost toho, že „chyba **h** je menší než  $\epsilon$  (**h** je  **$\epsilon$ -skoro správná**)“ je větší než  $(1 - \delta)$ .

### Praktické použití odhadu počtu příkladů

$m \geq 1/\epsilon * (\ln |H| + \ln (1/\delta))$

Necht'  $H_B(n)$  je množina všech bool. funkcí pro **n** bool. atributů, tj. zobrazení z **n**-tíc skládajících se z 0 a 1 do  $\{0, 1\}$ .

- Velikost definičního oboru:  $2^n$
- Počet funkcí z množiny mohutnosti **a** do  $\{0, 1\}$ :  $2^a$ .

Tedy mohutnost  $H_B(n)$  je  $2^{2^n}$

Postačující podmínka pro počet příkladů  $m_B(n)$ , které potřebujeme k tomu, abychom se skoro správně naučili koncept popsaný booleovskou funkcí o **n** atributech, je  
 $m_B(n) \geq 1/\epsilon * (2^n + \lg (1/\delta))$

**Důsledek:**  
 Máme-li se skoro správně naučit koncept popsaný obecnou bool. funkcí, pak potřebujeme více než  $2^n$  příkladů. Jinými slovy: musíme znát celý definiční obor. Věta o osklivém kačátku.

### Důsledek odhadu

Je-li **H** množina možných hypotéz, pak se lze skoro správně ( $\epsilon$  pravděpodobností větší než  $(1 - \delta)$ ) naučit hypotézu, jejíž chyba je menší než  $\epsilon$ , pokud máme **m** trénovacích příkladů a platí  
 $m \geq 1/\epsilon * (\ln |H| + \ln (1/\delta))$ . (i)

**Pozorování:** **m** je funkcí  $|H|$

Podají-li se nám získat nějakou doplňkovou informaci (omezení na tvar přípustných hypotéz), která omezuje rozsah **H**, pak vystačíme s menším počtem trénovacích příkladů !!! Zde hraje významnou roli doménová znalost.

Pokusme se

- provést odhad mohutnosti množiny hypotéz pro některé běžné typy hypotéz (rozhodovací stromy,...)
- a zjistit vliv tohoto odhadu na požadovaný počet trénovacích příkladů

### Věta o PAC-účení rozhodovacího stromu

Nechť objekty jsou charakterizovány pomocí  $n$  binárních atributů a nechť připouštíme jen hypotézy ve tvaru rozhodovacího stromu s maximální délkou větve  $k$ . Dále nechť  $\delta, \epsilon$  jsou malá pevně zvolená kladná čísla blízká 0. Pokud algoritmus strojového učení vygeneruje hypotézu  $\phi$ , která je konzistentní se všemi  $m$  příklady trénovací množiny a platí

$$m \geq m_{k-DT}(\delta) \geq c \left( n^k + \ln(1/\delta) \right) / \epsilon$$

pak  $\phi$  je  $\epsilon$ -skoro správná hypotéza s pravděpodobností větš než  $(1-\delta)$ , t.j. **chyba hypotézy  $\phi$  na celém definičním oboru konceptu je menší než  $\epsilon$  s pravděpodobností větš než  $(1-\delta)$ .**

*Doporučená literatura:*  
 Mařík et al.: **UI(3)** – kapitola 7 Teorie složitosti a úlohy UI (Demlová, Štěpánková)

### Odhad mohutnosti množiny hypotéz pro rozhodovací seznam (lin. reprezentace stromu)

$L_n$  jazyk obsahující přesně  $n$  binárních atributů  
 $\Omega$  def. obor uvažovaného konceptu má tedy  $2^n$  různých prvků

**Rozhodovací seznam** (decision list) v jazyce  $L_n$  je uspořádaný seznam  $\mathfrak{R} = [t_1:c_1, \dots, t_m:c_m]$ , kde

- $t_i$  je test vyjádřený ve tvaru konjunkce literálů z  $L_n$
- $c_i \in \{0,1\}$  je přiřazená klasifikace.

Nechť  $o \in \Omega$ , pak  $\mathfrak{R}(o) = c_i$ , kde  $t_i$  je první test, který objekt  $o$  splňuje (t.j.  $t_1(o)=0, \dots, t_{i-1}(o)=0, t_i(o)=1$ ).

Pokud  $t_k(o)=0$  pro všechna  $k \leq m$ , pak  $\mathfrak{R}(o) = 0$

*Každý strom hloubky  $n$  (= délka nejdelší větve) nebo formuli v disjunktivní normální formě lze napsat jako rozhodovací seznam v jazyce  $L_n$ ; např.*  
 $(s_1 \& s_3 \& \text{not } s_2)$  v  $(\text{not } s_1 \& s_3 \& \text{not } s_2)$  odpovídá  
 $[(s_1 \& s_3 \& \text{not } s_2):1, (\text{not } s_1 \& s_3 \& \text{not } s_2):1]$

### Odhad mohutnosti $k$ -DL( $n$ )

Nechť  $k$ -DL( $n$ ) je množina všech rozhodovacích seznamů, jejichž testy mají přípustnou délku omezenou pevně zvoleným číslem  $k < n$ . Jak ovlivní volba  $k$  mohutnost množiny hypotéz?

**Odhad mohutnosti prostoru hypotéz pro rozhodovací seznamy 1-DL( $n$ )**  
 $|1\text{-DL}(n)| < \text{počet permutací z } n \text{ (tedy } n!) \text{ krát } 3^n$ .

Zdůvodnění: Z 1 pevně zvolené permutace  $n$  prvků lze totiž vytvořit  $3^n$  různých rozhodovacích seznamů (test je použit s výsledkem 1, 0 nebo „nezařazen“).  
 $|1\text{-DL}(n)| < n! \cdot 3^n$

Protože  $\ln(n!) < n \cdot \ln n$ , platí  
 $\ln |1\text{-DL}(n)| < \ln(n! \cdot 3^n) < O(n \cdot \ln n) + n \ln 3$

Skoro správného učení lze dosáhnout při rozumném počtu trén. příkladů  
 $m \geq 1/\epsilon \cdot (\ln |H| + \ln(1/\delta))$ ; pro 1-DL( $n$ ) jde o číslo srovnatelné s  $(n \cdot \lg n)$ , což je výrazně méně než mohutnost  $2^n$  celého uvažovaného definičního oboru

### Odhad mohutnosti $k$ -DL( $n$ )

**Conj( $n, k$ )**: počet různých konjunkcí nejvýše  $k$  literálů sestavených z  $n$  atributů.  
**Conj0( $n, j$ )**: počet všech konjunkcí přesně  $j$  literálů sestavených z  $n$  atributů (pomocná veličina pro odhad **Conj( $n, k$ )**).

Postupujeme takto  
 $\text{Conj0}(n, j) < 2^j \cdot n^j = (2n)^j$  člen vyjadřující „znaménko“ atomu  
 $\text{Conj}(n, k) < \sum_{i=0}^k \text{Conj0}(n, i)$   
 $< \sum_{i=0}^k (2n)^i = 2n(2^k - 1) / (2n - 1) \approx O(n^k)$  (ii)

Horní odhad pro počet prvků  $k$ -DL( $n$ ): Rozhodovací seznam je vlastně uspořádaná posloupnost neopakujících se prvků z **Conj( $n, k$ )**, z nichž každý je klasifikován jednou z hodnot  $\{0,1, \times\}$ , kde „ $\times$ “ chápeme tak, že daná konjunkce v rozhodovacím seznamu není. Zřejmá tedy

$$|k\text{-DL}(n)| < 3 \cdot |\text{Conj}(j, k)| \cdot |\text{Conj}(j, k)|$$

### Odhady mohutnosti $k$ -DL( $n$ ) a $|k$ -DL( $n$ )

Víme, že  $|k\text{-DL}(n)| < 3 \cdot |\text{Conj}(j, k)| \cdot |\text{Conj}(j, k)|$ . Z toho plyne, že  
 $\ln |k\text{-DL}(n)| < |\text{Conj}(j, k)| \cdot \ln 3 + \ln |\text{Conj}(j, k)|$

Použitím vztahu  $\lg n! < n \cdot \lg n$  dostáváme  
 $\ln |k\text{-DL}(n)| < |\text{Conj}(j, k)| \cdot (\ln 3 + \ln |\text{Conj}(j, k)|)$   
 $< O(n^k) \cdot (\ln 3 + \ln O(n^k)) \approx O(n^k \ln(n^k))$

Po dosazení do vzorce  $m \geq 1/\epsilon \cdot (\ln |H| + \ln(1/\delta))$  dostáváme odhad pro hypotézy ve tvaru rozhodovacího listu  
 $m_{k-DL}(\delta) \geq c/\epsilon \cdot (O(n^k \ln n^k) + (1/\delta))$

Pro rozhodovací stromy s omezenou hloubkou je odhad ještě poněkud nižší, protože pro mohutnost prostoru hypotéz platí  
 $|k\text{-DT}(n)| < 3 \cdot |\text{Conj}(j, k)|$   
 Odpovídající počet trén. příkladů je  $m_{k-DT}(\delta) \geq c/\epsilon \cdot (n^k + (1/\delta))$