

FAKULTA ELEKTROTECHNIKY A KOMUNIKAČNÍCH TECHNOLOGIÍ
VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ

BAYESOVSKÉ UČENÍ

Autor textu:
Ing. Petr Honzík, Ph.D.

Květen 2014

Komplexní inovace studijních programů a zvyšování kvality výuky na FEKT VUT v Brně
OPVK CZ.1.07/2.2.00/28.0193



INVESTICE DO ROZVOJE VZDĚLÁVÁNÍ

Obsah přednášky

- 1. Bayesův teorém**
2. Brutální Bayesovský klasifikátor (BBK)
3. Maximální aposteriorní pravděpodobnost (MAP)
4. Optimální Bayesovský klasifikátor (OBK)
5. Gibbsův algoritmus (GiA)
- 6. Naivní Bayesovský klasifikátor (NBK)**

Thomas Bayes

- duchovní, 2. polovina 18. století
- zavrhován statistiky (nepodložené empirické metody)
- zabýval se otázkou, jak pozdější zkušenosti uvést do souladu s původními předpoklady (dynamické „ověřování hypotéz“ s jejich dodatečnou korekcí)
- **(a)priorní pravděpodobnost** (dána na počátku) vs. **(a)posteriovní pravděpodobnost** (vyplývající z následné analýzy, ověřené zkušeností)
- Principiálně blízké přirozenému lidskému uvažování

Bayesův vzorec

$$p(H_k | D) = \frac{p(D | H_k) \cdot p(H_k)}{\sum_{i=1}^K p(D | H_i) \cdot p(H_i)} = p(H_k) \frac{p(D | H_k)}{p(D)}$$

H_k je **hypotéza**

D je realita, data, konkrétní měření

$p(H_k)$ zkušenost; **apriorní pravděpodobnost**

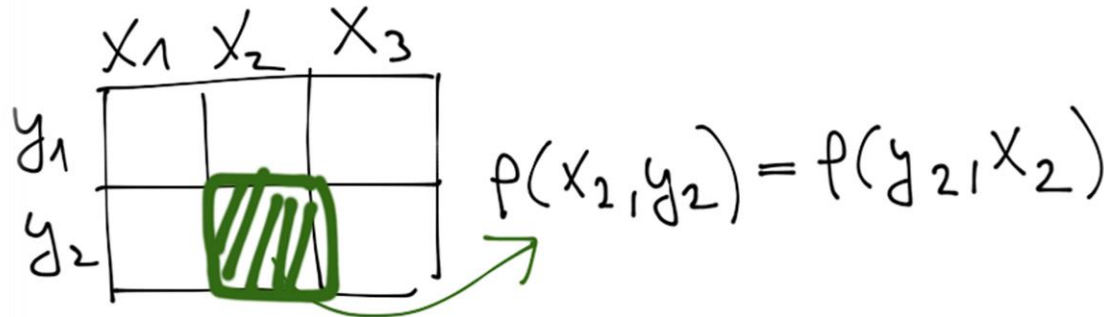
$p(D | H_k)$ je **věrohodnost**

$p(D)$ pravděpodobnost nastolení měření (dat, údajů) D , **evidence**

$p(H_k | D)$ je **aposteriorní pravděpodobnost** platnosti konkrétní hypotézy

? uveď a popiš Bayesův vzorec

Bayesův vzorec ze sdružené pravděpodobnosti

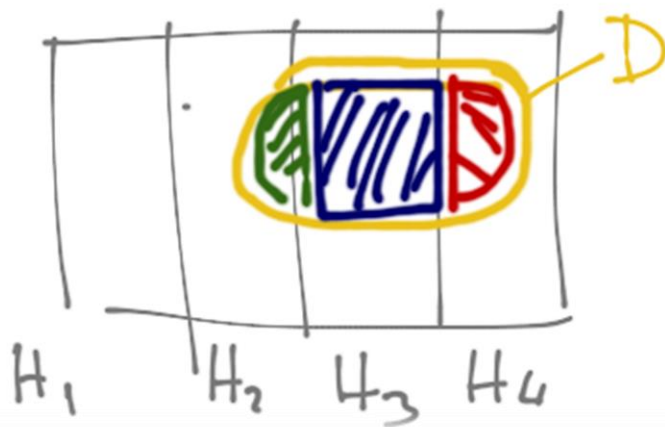


$$p(x_2, y_2) = p(x_2 | y_2) \cdot p(y_2)$$

$$p(y_2, x_2) = p(y_2 | x_2) \cdot p(x_2)$$

$$\Rightarrow p(x_2 | y_2) = \frac{p(y_2 | x_2) \cdot p(x_2)}{p(y_2)}$$

Bayesův vzorec graficky – 1/3

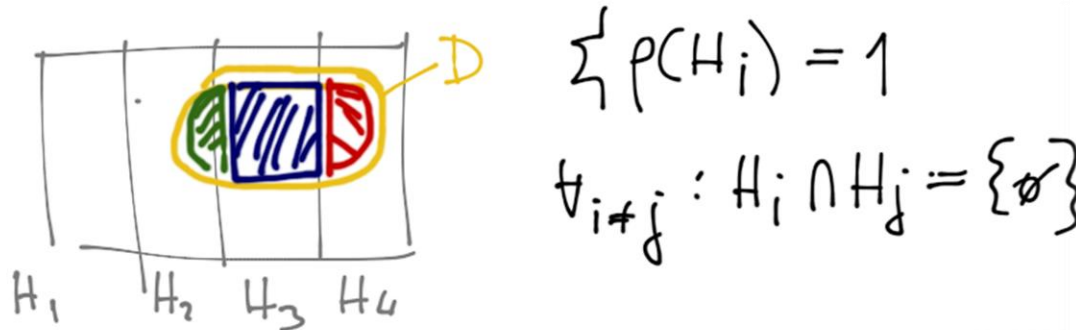


$$\sum_i P(H_i) = 1$$

$$\forall i \neq j : H_i \cap H_j = \{\emptyset\}$$

? uveď a popiš Bayesův vzorec

Bayesův vzorec graficky – 2/3



 = $p(D)$

 = $p(D|H_2) \cdot p(H_2)$

$$\Rightarrow p(H_2|D) = \frac{\text{green oval}}{\text{yellow oval}} = \frac{p(D|H_2) p(H_2)}{p(D)}$$

? uveď a popiš Bayesův vzorec

Bayesův vzorec graficky – 3/3



$$P(D) = P(D|H_2) \cdot P(H_2) + P(D|H_3) \cdot P(H_3) + P(D|H_4) \cdot P(H_4)$$

$$\Rightarrow P(H_2|D) = \frac{P(D|H_2) \cdot P(H_2)}{\sum_i P(D|H_i) \cdot P(H_i)}$$

? uveď a popiš Bayesův vzorec

Bayesovské učení - poznámky

- Klademe si dvě základní otázky:
 - Jaká **hypotéza** (model) o modelovaném systému je s největší pravděpodobností platná?
 - Jaká je **predikce** nové instance na základě známých hypotéz (modelů)?
- Výhody
 - lze kombinovat **předchozí znalost** s **pozorovanými měřeními** (apriorní pravděpodobnost, evidence)
 - intuitivní řešení **blízké lidskému uvažování**
- Nevýhody
 - narůstající **složitost** hledání řešení s rostoucím počtem hypotéz – v obecném případě

Příklad 1 – Bayes

- V pytlíku je 7 hracích kostek. 4 jsou normální, na 1 padají víc šestky, na 2 víc jedničky.
- Jaká je pravděpodobnost, že vytáhnu normální kostku? (apriorní pravděpodobnost)
- **Jaká je pravděpodobnost, že když po 20 hodech (experiment E) padla šestka 6-krát, mám kostku *Normální, Šestkovou, Jedničkovou?*** (aposteriorní pravděpodobnost)
- (když $p(6/N)=1/6$, $p(6|\check{S})=1/3$, $p(6|J)=2/15$)

$$p(N/E)=0,52; p(\check{S}/E)=0,36; p(J/E)=0,12$$

Ale co teď s tím?

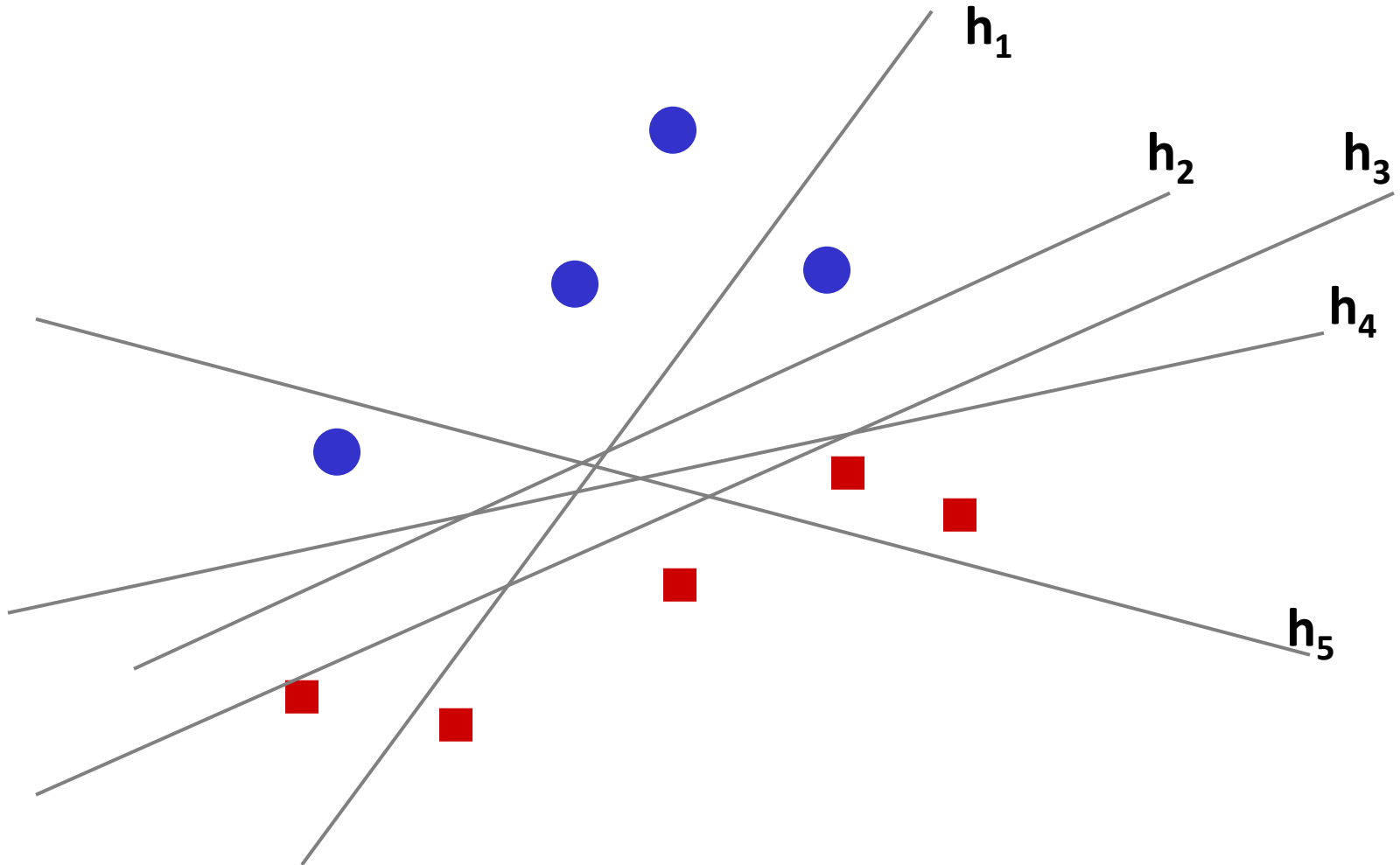
- Víme, že s
 - $p=0,52$ držíme v ruce kostku normální
 - $p=0,36$ kostku šestkovou
 - $p=0,12$ kostku jedničkovou.
- Jak z toho ale predikovat?
- Jak odhadnout pravděpodobnost, s jakou padne číslo „6“ v dalším hodu?

Brutální Bayesovský klasifikátor

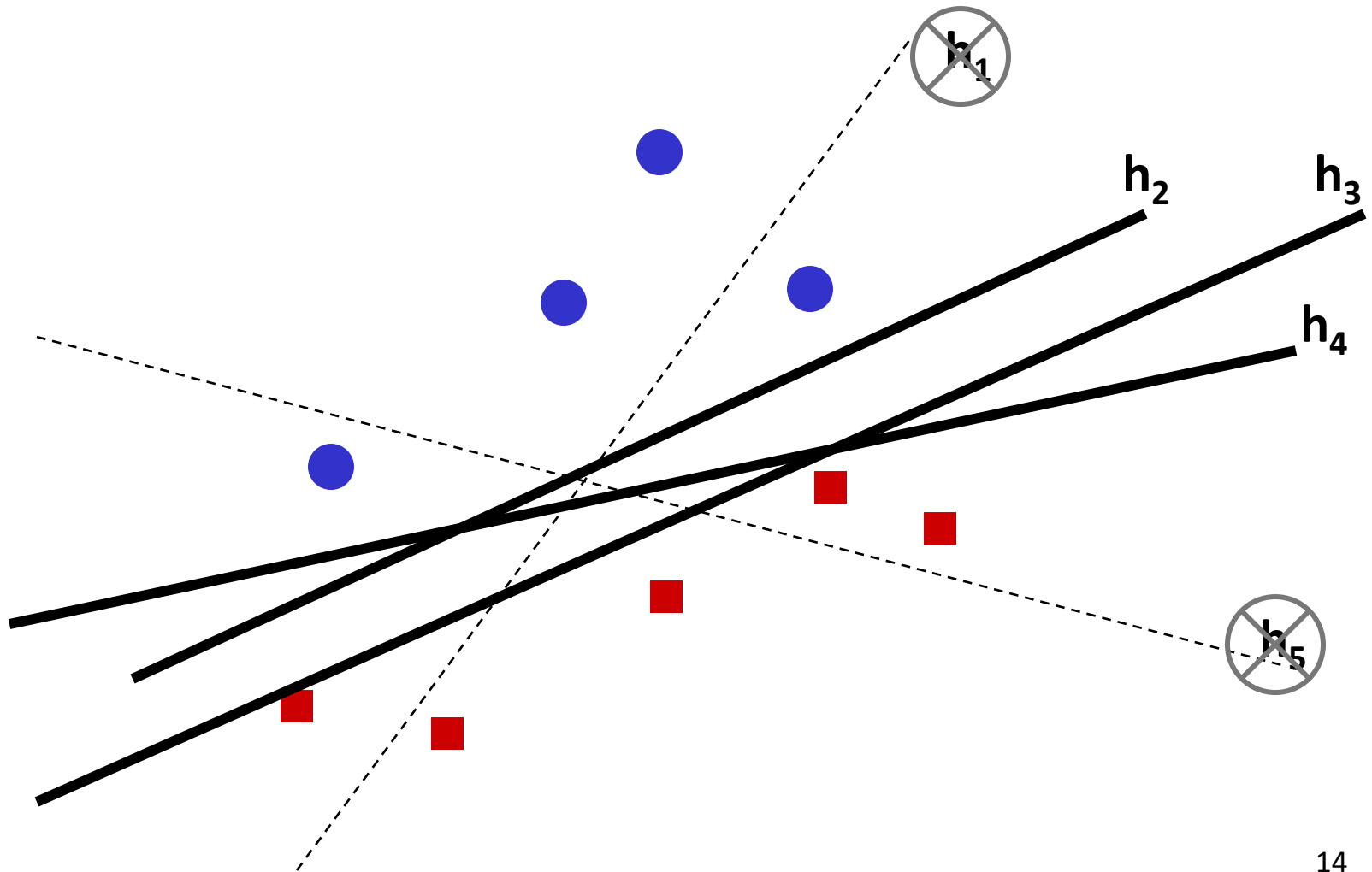
- Princip
 - **Všechny hypotézy konzistentní s daty mají stejnou pravděpodobnost** $P(h_j)$, hypotézy nekonzistentní jsou vyloučeny
 - Pokud jsou trénovací data zašuměná, budou správné hypotézy zamítnuty, možné zamítnutí všech hypotéz
- Predikce
 - pokud K modelů predikuje na trénovacích datech správně, jsou si podle BBK tyto hypotézy rovnocenné; predikce je průměr nebo nejčastější třída

$$\left. \begin{aligned}
 P(D|h_j) &= \begin{cases} 1; \forall y_i : y_i = h_j(x_i) \\ 0; \exists y_i : y_i \neq h_j(x_i) \end{cases} \\
 P(h_j) &= \frac{1}{\sum_{j=1}^{\text{hypotez}} P(D|h_j)} = \frac{1}{|H_{konz}|} = konst
 \end{aligned} \right\} \Rightarrow g_{out} = \arg \max_{j=1 \dots C} \sum_{h_i \in H} P(g_j|h_i)$$

BBK – příklad



BBK – příklad



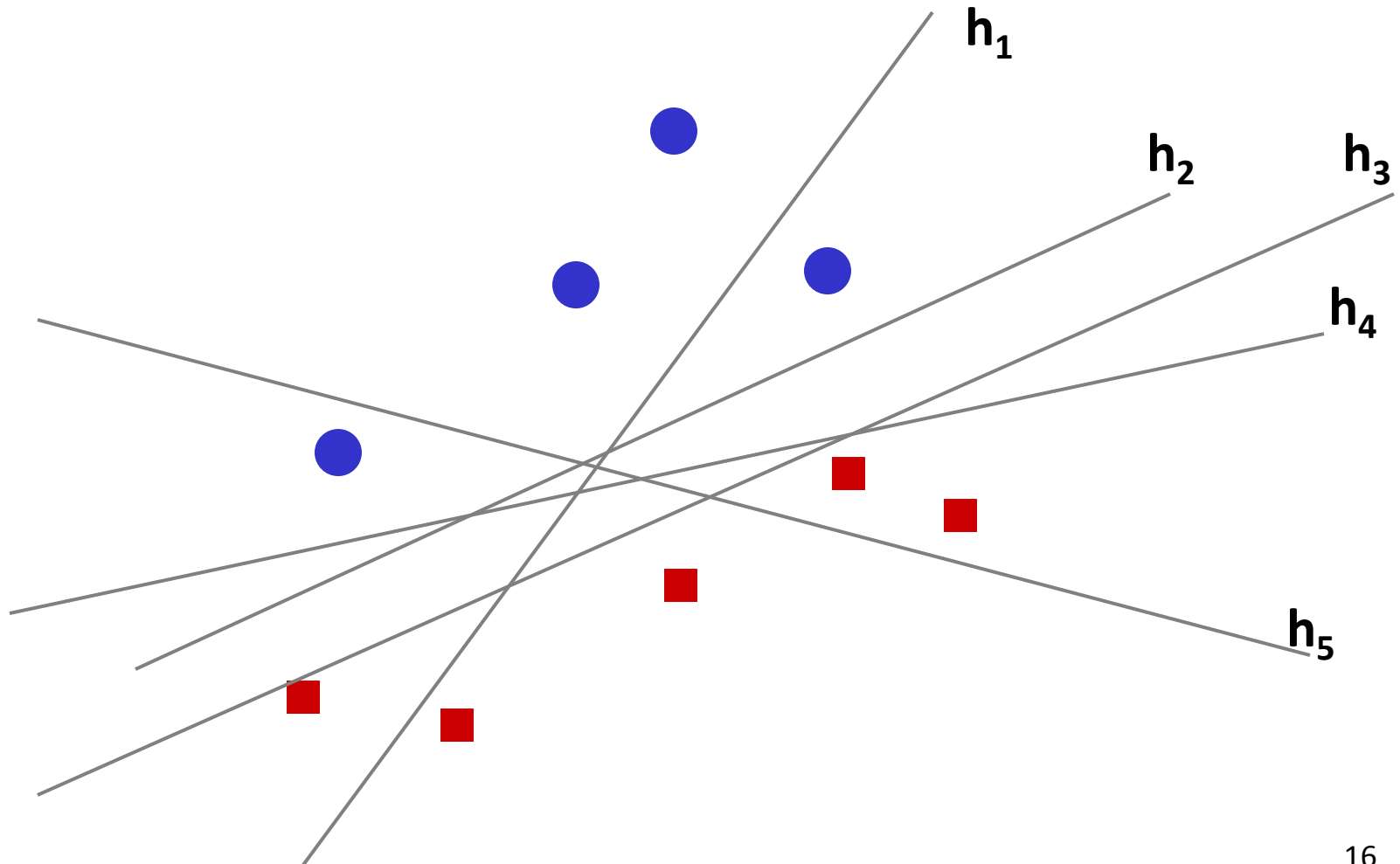
Maximální a posteriorní pravděpodobnost

- **Nejpravděpodobněji platná hypotéza** je ta s největší a posteriorní pravděpodobností.
- Podoba s maximální věrohodností a MNČ

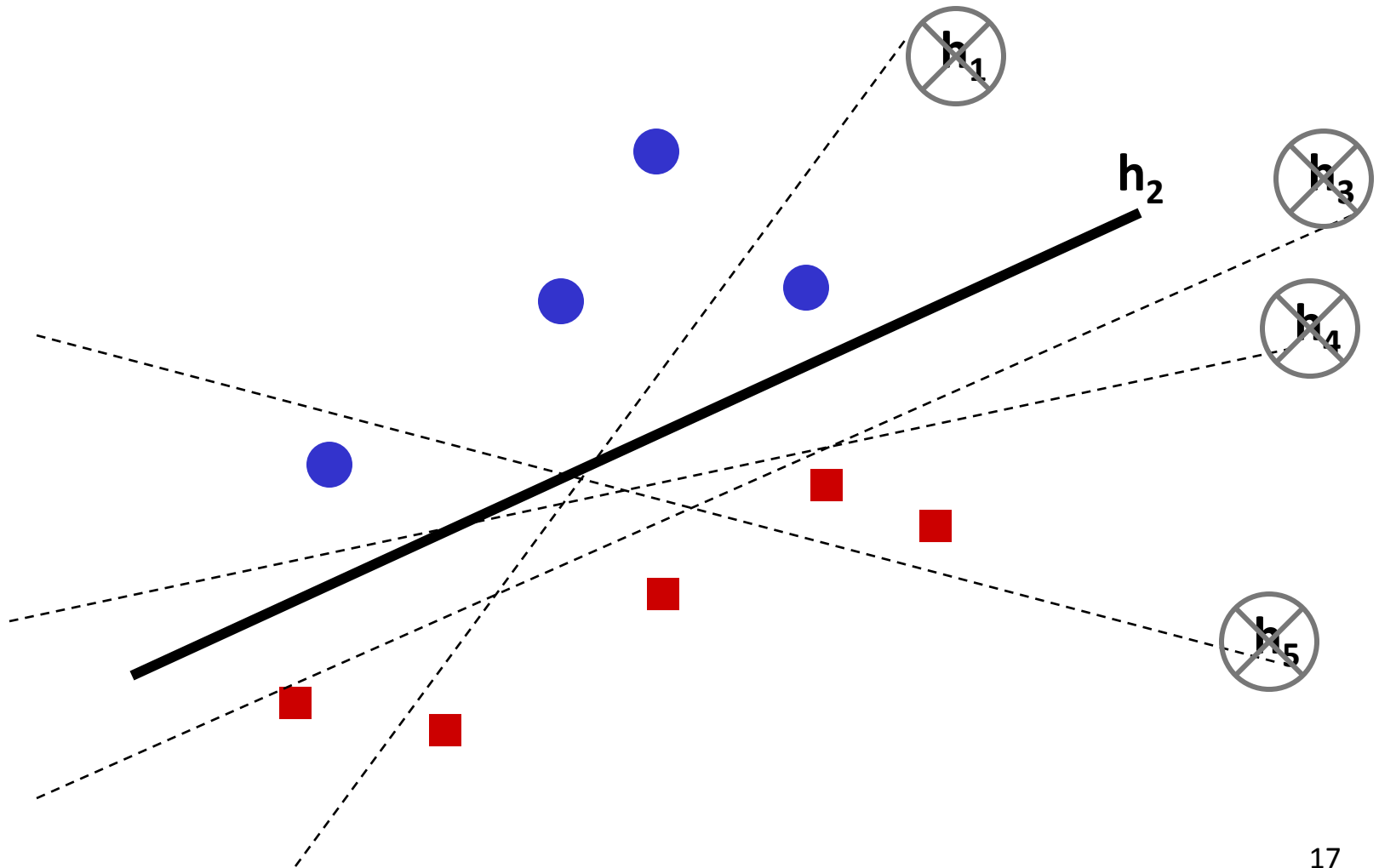
$$\begin{aligned}h_{MAP} &= \arg \max_{j=1 \dots |H|} P(h_j | D) = \\ &= \arg \max_{j=1 \dots |H|} P(h_j) \frac{P(D | h_j)}{P(D)} = \\ &= \arg \max_{j=1 \dots |H|} P(h_j) \cdot P(D | h_j)\end{aligned}$$

$$g_{out} = \arg \max_{j=1 \dots C} P(g_j | h_{MAP})$$

MAP – příklad



MAP – příklad

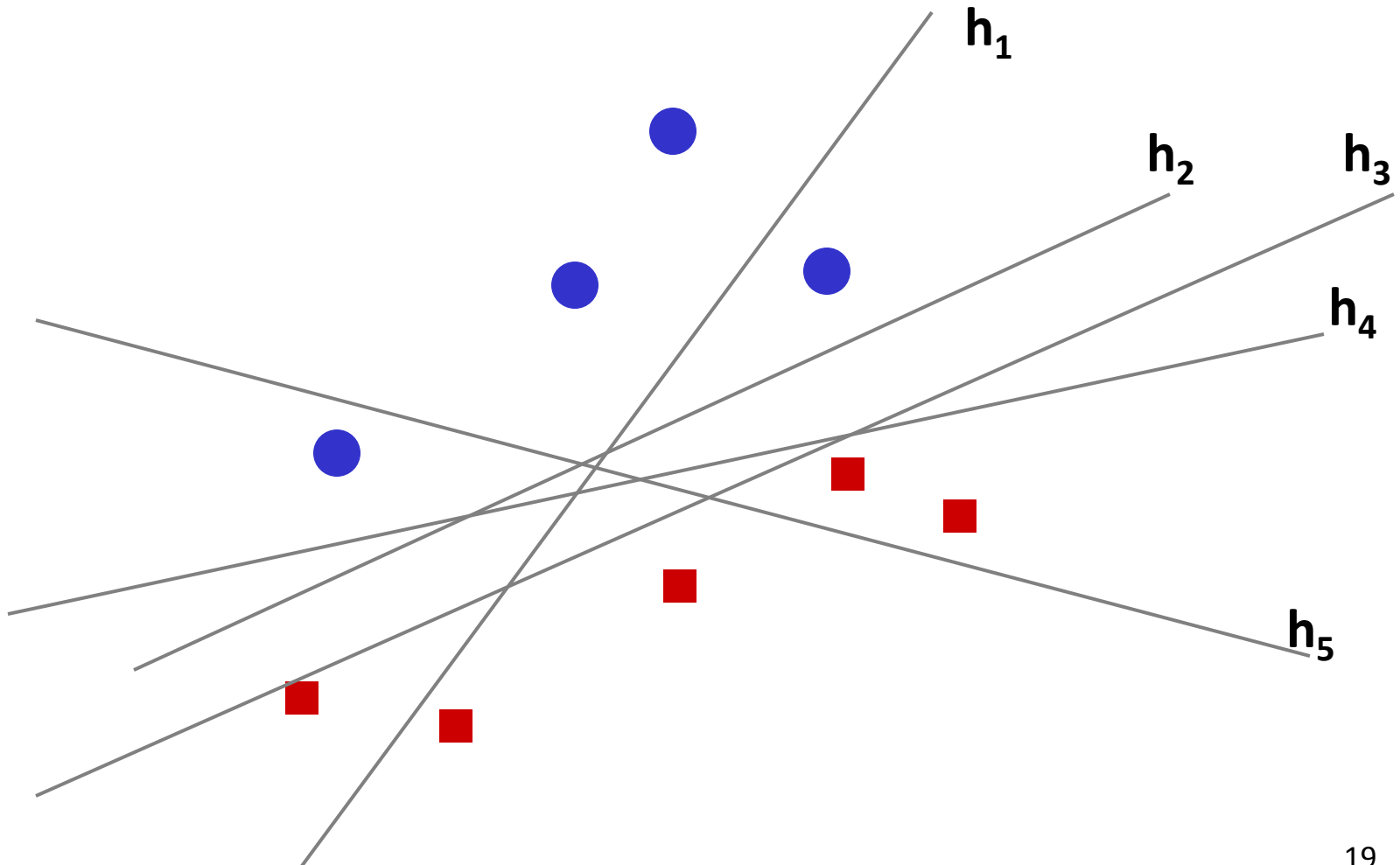


Optimální Bayesovský klasifikátor

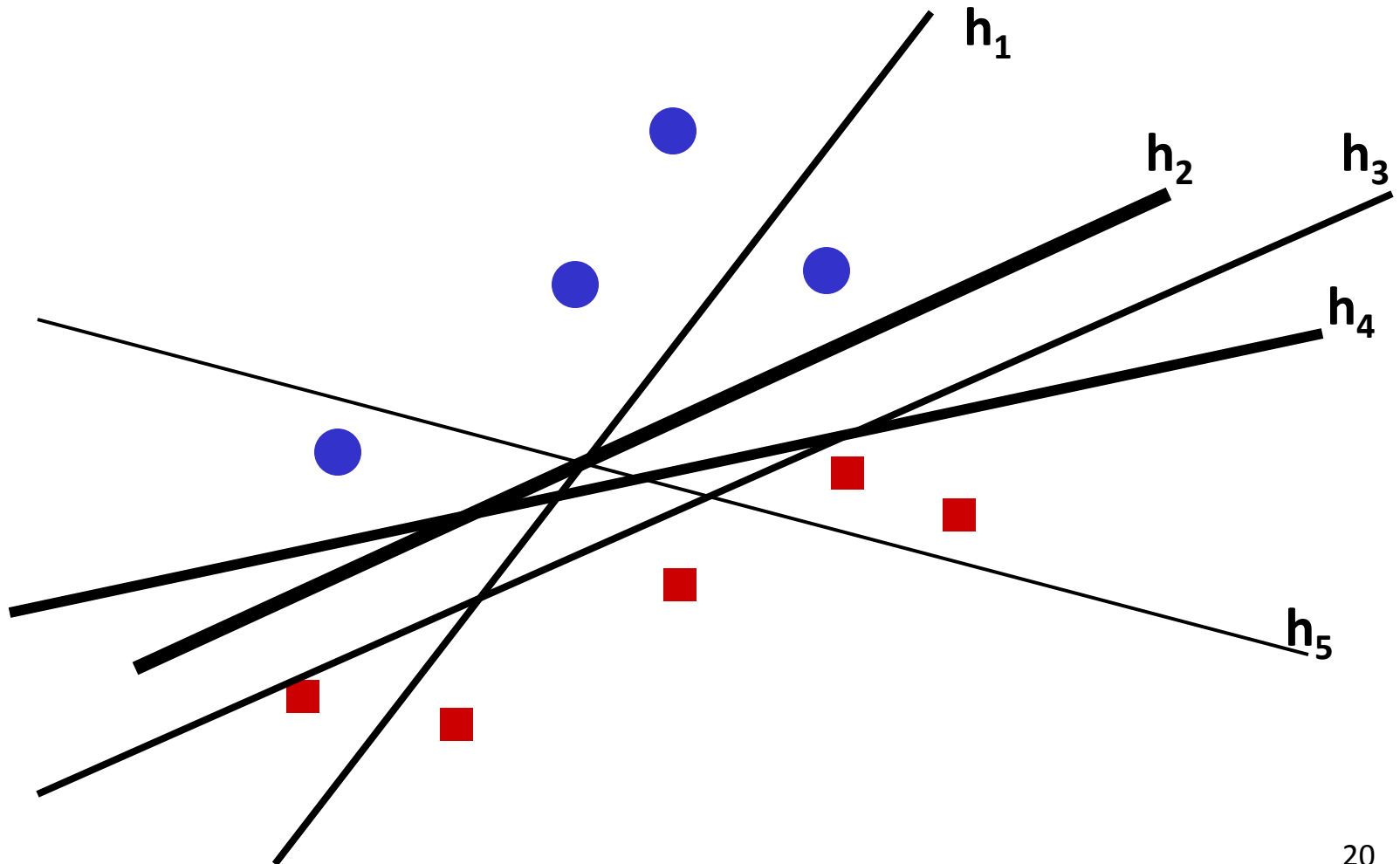
- „Jaká je **nejpravděpodobnější klasifikace g_j** nové instance za předpokladu informací založených na několika hypotézách h_i ?“
- jinými slovy, optimální bayesovský klasifikátor vychází z hypotéz a jejich věrohodností; výpočty sejné jako u MAP, nicméně není vybrán 1 nejlepší, ale **všechny hypotézy se úměrně své kvalitě podílejí na predikci**
- Klasifikujme binárně do tříd +, – a mějme hypotézy s aposteriorní pravděpodobností $p(h_1 | D)=0,4$ $p(h_2 | D)=0,3$ a $p(h_3 | D)=0,3$. Pokud podle h_1 je výsledek negativní a podle h_2 a h_3 je výsledek pozitivní, MAP považuje za výsledek predikci dle h_1 , OBK však dle h_2 a h_3 , protože $(0,3+0,3) > 0,4$.

$$g_{\max} = \arg \max_{j=1..C} \sum_{h_i \in H} P(g_j | h_i) \cdot P(h_i | D)$$

OBK – příklad



OBK – příklad



Gibbsův algoritmus

- řeší hlavní problém OBK – tím je **výpočetní náročnost**
 - použití všech váhovaných modelů
- Princip
 - při každé predikci vyber náhodně (avšak úměrně posteriorní pravděpodobnosti modelů) jednu hypotézu a podle té predikuj
 - Dosažená chyba je maximálně dvojnásobná oproti OBK

Příklad 2 – BBK, MAP, OBK

- V pytlíku je 7 hracích kostek. 4 jsou normální, na 1 padají víc „6“, na 2 víc „1“.
- Jaká je pravděpodobnost, že když po **experimentu E** ve 20 hodech padla šestka 6-krát, mám kostku *Normální, Šestkovou, Jedničkovou*? (aposteriorní pravděpodobnost)
- (když $p(6|N)=1/6$, $p(6|\check{S})=1/3$, $p(6|J)=2/15$); $p(N|E)=0,52$ $p(\check{S}|E)=0,36$ $p(J|E)=0,12$
- **S jakou pravděpodobností padne šestka ve 21. hodů?**

- **BBK:** • Konzistentní (možné) jsou všechny hypotézy,
 $p=(1/6+1/3+2/15)/3=0,21$
- **MAP:** • Nejpravděpodobnější byla hypotéza „normální“,
tedy $1/6$, $p=0,17$
- **OBK:** • Výsledná pravděpodobnost je dána součtem
součinů pravděpodobnosti hypotézy a nastolení
dané události, tedy
 $p = 0,52 \cdot 1/6 + 0,36 \cdot 1/3 + 0,12 \cdot 2/15 = 0,22$

Příklad 3

V pytlíku mohou být jehlany j a kostičky k . Vytáhli jsme 1 kostičku – **experiment E**. S jakou pravděpodobností p_k vytáhneme další kostičku?

- Víme, že v pytlíku bylo na počátku právě 5 objektů
- Víme, že na počátku je jakákoliv kombinace jehlanů a kostiček v pytlíku stejně pravděpodobná

h_i = hypotéza, že v pytlíku bylo i -kostiček;

urči všechny $p(h_i|E)$ a pro jednotlivé metody p_k v dalším tahu

- **BBK:** Až na h_0 jsou konzistentní; $p(k|h_1)=0$; $p(k|h_2)=0,25$...

$$p_k = \sum p(k|h_i)/5 = (0+0,25+0,5+0,75+1)/5 = 0,5$$
- **MAP:** Nejpravděpodobnější je h_5 , tedy $p_k=1$
- **OBK:** $p(E|h_0)=0$; $p(E|h_1)=0,2$; ... $p(E|h_5)=1$;
 protože $p(h_i)=\text{konst}$, platí $p(h_i|E) = \frac{p(E|h_i) \cdot p(h_i)}{\sum p(E|h_i) \cdot p(h_i)} = \frac{p(E|h_i)}{\sum p(E|h_i)}$

$$p_k = \sum [p(k|h_i) p(E|h_i) / \sum p(E|h_i)] = 0,67$$

Naivní Bayesovský klasifikátor - NBK

- navzdory předpokladu nezávislosti veličin (není většinou pravda) velice přesný; kvůli zjednodušení slovo *naivní*
- vychází z Bayesovy podmíněné pravděpodobnosti
- **klasifikace** na základě nejpravděpodobnější klasifikace g_{MAP} **za předpokladu vstupního vektoru** (x_1, \dots, x_n)

$$g_{MAP} = \arg \max_{j=1 \dots C} P(g_j | x_1, \dots, x_n)$$

C je počet tříd

n je počet atributů

Proč se NBK jmenuje naivní? V čem naivita spočívá?

Naivní Bayesovský klasifikátor - NBK

- což lze na základě bayesovy věty upravit následujícím způsobem:

$$\begin{aligned}g_{MAP} &= \arg \max_{j=1\dots C} P(g_j) \frac{P(x_1, \dots, x_n | g_j)}{P(x_1, \dots, x_n)} = \\ &= \arg \max_{j=1\dots C} P(g_j) \cdot P(x_1, \dots, x_n | g_j)\end{aligned}$$

- $P(g_j)$ se určí na základě četnosti výskytu v trénovacích datech
- $P(x_1, \dots, x_n | g_j)$ pro větší n (desítky) prakticky nemožné zjistit (tolik dat nemáme...). Určí se proto na základě zjednodušeného předpokladu, že hodnoty vstupních veličin jsou na sobě podmíněně nezávislé (což není většinou pravda)

Naivní Bayesovský klasifikátor - NBK

- pro výpočet sdružené podmíněné pravděpodobnosti platí:

$$\begin{aligned} P(x_1, \dots, x_n | g_j) &= P(x_1 | g_j) \cdot P(x_2, \dots, x_n | g_j, x_1) = \dots = \\ &= P(x_1 | g_j) \cdot P(x_2 | g_j, x_1) \cdot P(x_3 | g_j, x_1, x_2) \dots \end{aligned}$$

- to lze za předpokladu nezávislosti **zjednodušit** takto:

$$P(x_1, \dots, x_n | g_j) = \prod_{i=1..n} P(x_i | g_j)$$

- konečný vztah pro výslednou klasifikaci je dán vztahem

$$g_{MAP} = \arg \max_{j=1..C} P(g_j) \cdot \prod_{i=1..n} P(x_i | g_j)$$

? uveď rozdíl mezi optimálním a naivním Bayesovským klasifikátorem

NBK – okrajové podmínky

- **Problém**: situace, kdy konkrétní hodnota atributu pro danou třídu nikdy nenastala (žádná bruneta se nespálila); celý součin je pak kvůli jednomu členu roven 0, což není vhodné;

$$P(x_2|g_1) = 0 \Rightarrow P(x_1, \dots, x_n | g_1) = \prod_{i=1..n} P(x_i | g_1) = 0$$

- **Řešení**: výpočet podmíněné pravděpodobnosti dán vztahem:

$$P(x_1 | g_1) = \frac{n_c + mp}{n + m} \quad \text{tedy} : \frac{n_c + 1}{n + |\text{slovník}|}$$

kde n je počet prvků třídy g_1 (počet spálených), n_c počet prvků v g_1 s danou hodnotou atributu x_1 (spálených brunet), m váha (např. 1) a $p=1/k$, kde k je výčet hodnot atributu x_1 (3: bruneta, blond, zrzavá)

NBK – příklad

den	předpověď	teplota	vlhkost	vítr	hrát tenis?
1.	slunečno	teplo	vysoká	slabý	NE
2.	slunečno	teplo	vysoká	silný	NE
3.	zataženo	teplo	vysoká	slabý	ANO
4.	déšť	středně	vysoká	slabý	ANO
5.	déšť	chladno	normální	slabý	ANO
6.	déšť	chladno	normální	silný	NE
7.	zataženo	chladno	normální	silný	ANO
8.	slunečno	středně	vysoká	slabý	NE
9.	slunečno	chladno	normální	slabý	ANO
10.	déšť	středně	normální	slabý	ANO
11.	slunečno	středně	normální	silný	ANO
12.	zataženo	středně	vysoká	silný	ANO
13.	zataženo	teplo	normální	slabý	ANO
14.	déšť	středně	vysoká	silný	NE

$$g_{MAP} = \arg \max_{j=1..C} P(g_j) \cdot \prod_{i=1..n} P(x_i | g_j)$$

Bude se hrát 15. den? <slunečno,chladno,vysoká,silný>

NBK – příklad

$$P(\text{ANO})=9/14$$

$$P(\text{slunečno}/\text{ANO})=2/9$$

$$P(\text{chladno}/\text{ANO})=3/9$$

$$P(\text{vysoká}/\text{ANO})=3/9$$

$$P(\text{silný}/\text{ANO})=3/9$$

$$P(\text{NE})=5/14$$

$$P(\text{slunečno}/\text{NE})=3/5$$

$$P(\text{chladno}/\text{NE})=1/5$$

$$P(\text{vysoká}/\text{NE})=4/5$$

$$P(\text{silný}/\text{NE})=3/5$$

$$P(\text{ANO}) \cdot P(\text{slunečno}/\text{ANO}) \cdot P(\text{chladno}/\text{ANO}) \cdot P(\text{vysoká}/\text{ANO}) \cdot P(\text{silný}/\text{ANO}) = 0,0053$$

$$P(\text{NE}) \cdot P(\text{slunečno}/\text{NE}) \cdot P(\text{chladno}/\text{NE}) \cdot P(\text{vysoká}/\text{NE}) \cdot P(\text{silný}/\text{NE}) = 0,0206$$

$$P(\text{ANO} | \text{slunčeno, chladno, vysoká, silný}) = 0,0053 / (0,0053 + 0,0206) = 0,20$$

$$P(\text{NE} | \text{slunčeno, chladno, vysoká, silný}) = 0,0206 / (0,0053 + 0,0206) = 0,80$$

Princip jednoduchého spamového filtru

<i>z databáze slov</i>	spam (1000)	ham(300)
...
<i>mamka</i>	1	12
<i>ahoj</i>	60	23
<i>oběd</i>	2	3

- **Apriorní P:**

- $P(\text{spam}) = 0,9$

- $P(\text{ham}) = 0,1$

- $P(\text{spam} | \text{mamka,ahoj,oběd} - \text{bez normalizace}) = P(\text{spam}) * P(\text{mamka} | \text{spam}) * P(\text{ahoj} | \text{spam}) \dots$

- $P(\text{spam} | \dots - \text{bez normalizace}) = 0,9 * 1/1000 * 60/1000 * 2/1000 = 1,08E-7$

- $P(\text{ham} | \dots - \text{bez normalizace}) = 0,1 * 12/300 * 23/300 * 3/300 = 3,07E-6$

- **$P(\text{spam} | \text{mamka,ahoj,oběd}) =$**

$$P(\text{spam} | \text{mamka,ahoj,oběd}) / [P(\text{spam} | \text{mamka,ahoj,oběd}) + P(\text{ham} | \text{mamka,ahoj,oběd})] = 1,08E-7 / (1,08E-7 + 3,07E-6) = \mathbf{0,034}$$

- **$P(\text{ham} | \text{mamka,ahoj,oběd}) = 0,966$**

Shrnutí Bayese

- MAP (maximální aposteriorní pravděpodobnost) **vybír**á **nejpravděpodobnější hypotézu**. Podle této **jediné hypotézy** je následně predikováno.
- Brutální Bayesovské učení konceptů **vybír**á z množiny hypotéz (modelů) **hypotézy konzistentní** s trénovacími daty, které jsou si **rovnocenné**. Při predikci jsou si i **výstupy hypotéz rovnocenné**.
- Optimální Bayesovský klasifikátor stanoví **posterio**rní **pravděpodobnost všech hypotéz**. Při predikci jsou výstupy hypotéz **váhovány** posteriorními pravděpodobnostmi.
- Naivní Bayesovský klasifikátor **vytv**áří **klasifikátor**, predikuje na základě nové instance a trénovacích dat.