

Internet a klasifikační metody, 3. přednáška

Hlavní typy klasifikačních metod

volitelný předmět pro magisterské studium

Martin Holeňa



O čem to bude?

- ◆ Základní rozdělení klasifikačních metod do skupin
- ◆ Klasifikační metody nehledající hranici mezi třídami
 - klasifikace na základě k nejbližších sousedů
 - odhadování pravděpodobnosti tříd: bayesovská klasifikace
 - prokládání normálních rozdělání: diskriminační analýza
- ◆ Příklad metody hledající hranice: neuronové sítě

Rozdělení klasifikačních metod

◆ Kriterium: hledá metoda hranici mezi třídami?

1. *Primárně nehledají hranici*, ale zkušenost podobných

případů | pravděpodobnost tříd (bodové odhady, rozdělení)

- tradiční metody (předpočítačové), výpočetně nepříliš náročné

2. *Primárně hledají hranici*: moderní, výpočetně náročné

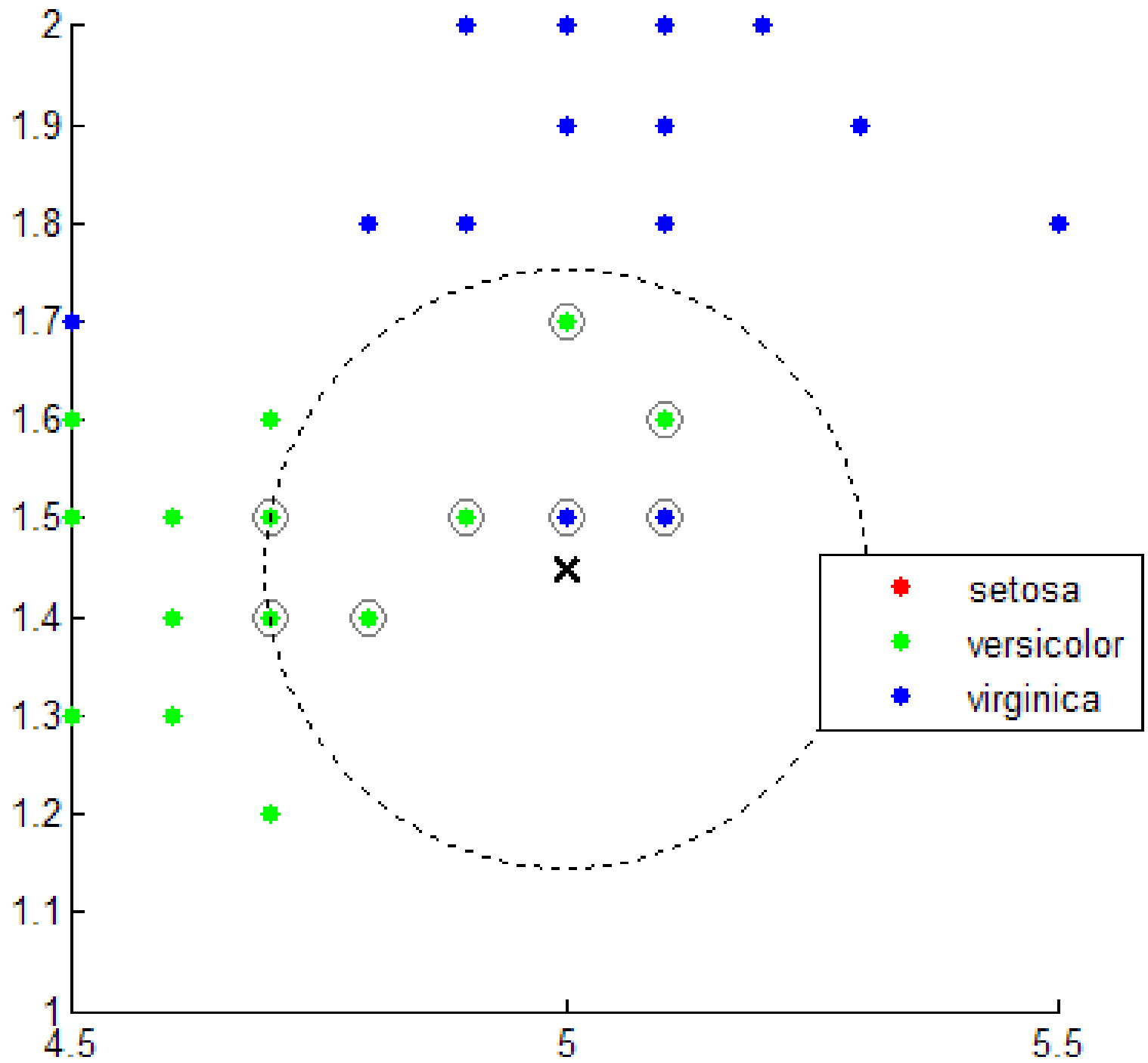
- např. neuronové sítě, rozhodovací stromy, SVM, ...

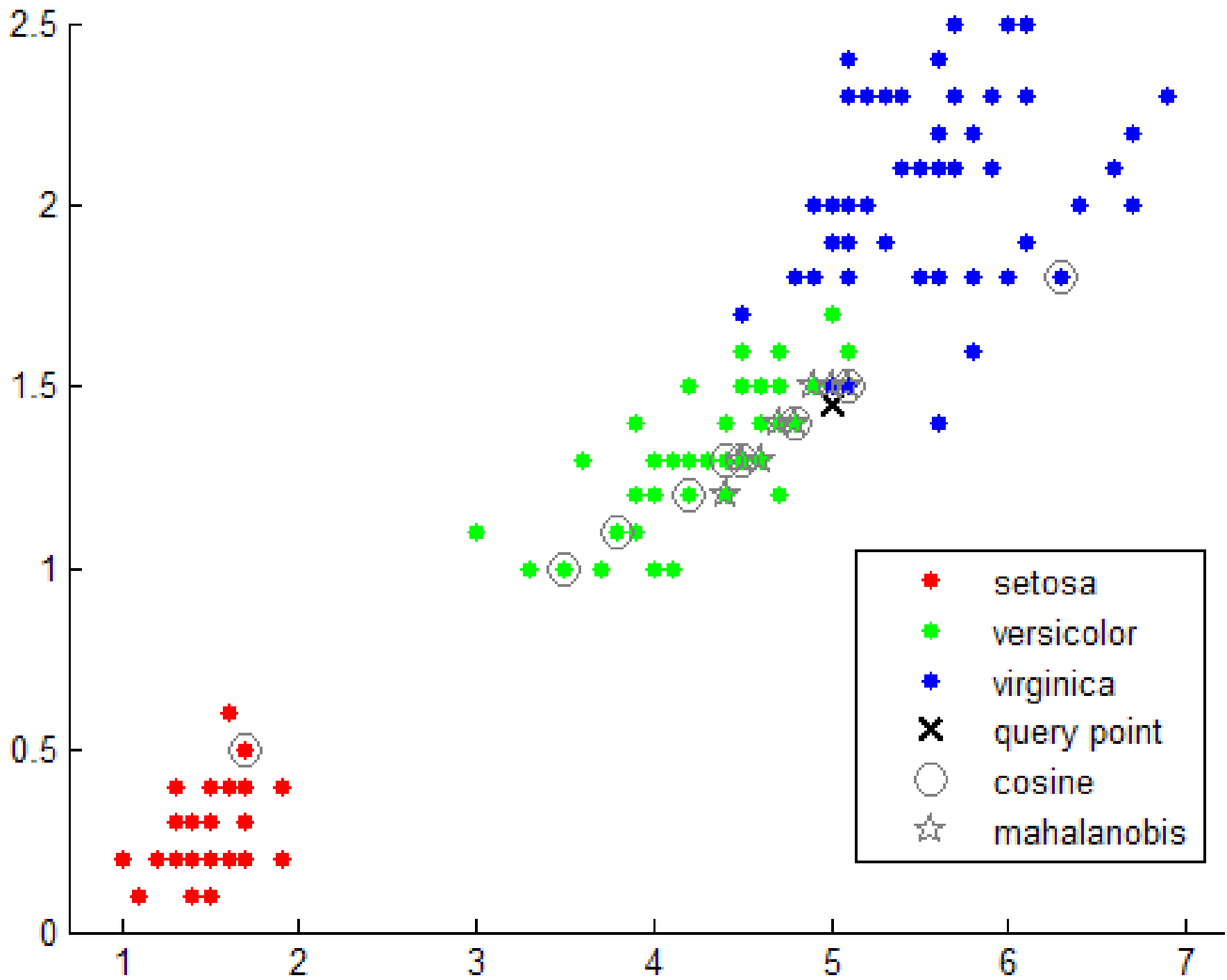
Klasifikace dle podobných případů

◆ Nový objekt klasifikován jako *třída nejčastější*

u k objektů nejpodobnějších podle zvolené podobnosti

- podmnožina objektů použitých pro učení (= zapamatování)
- \Rightarrow metoda k nejbližších sousedů (k nearest neighbours, k -NN)





Klasifikace dle podobných případů

- ◆ Nový objekt klasifikován jako *třída nejčastější*

u k objektů nejpodobnějších podle zvolené podobnosti

- podmnožina objektů použitých pro učení (= zapamatování)
- \Rightarrow metoda k nejbližších sousedů (k nearest neighbours, k -NN)

- ◆ Často silná závislost na volbě k

- *volba k* z konečného počtu možností: krosvalidací

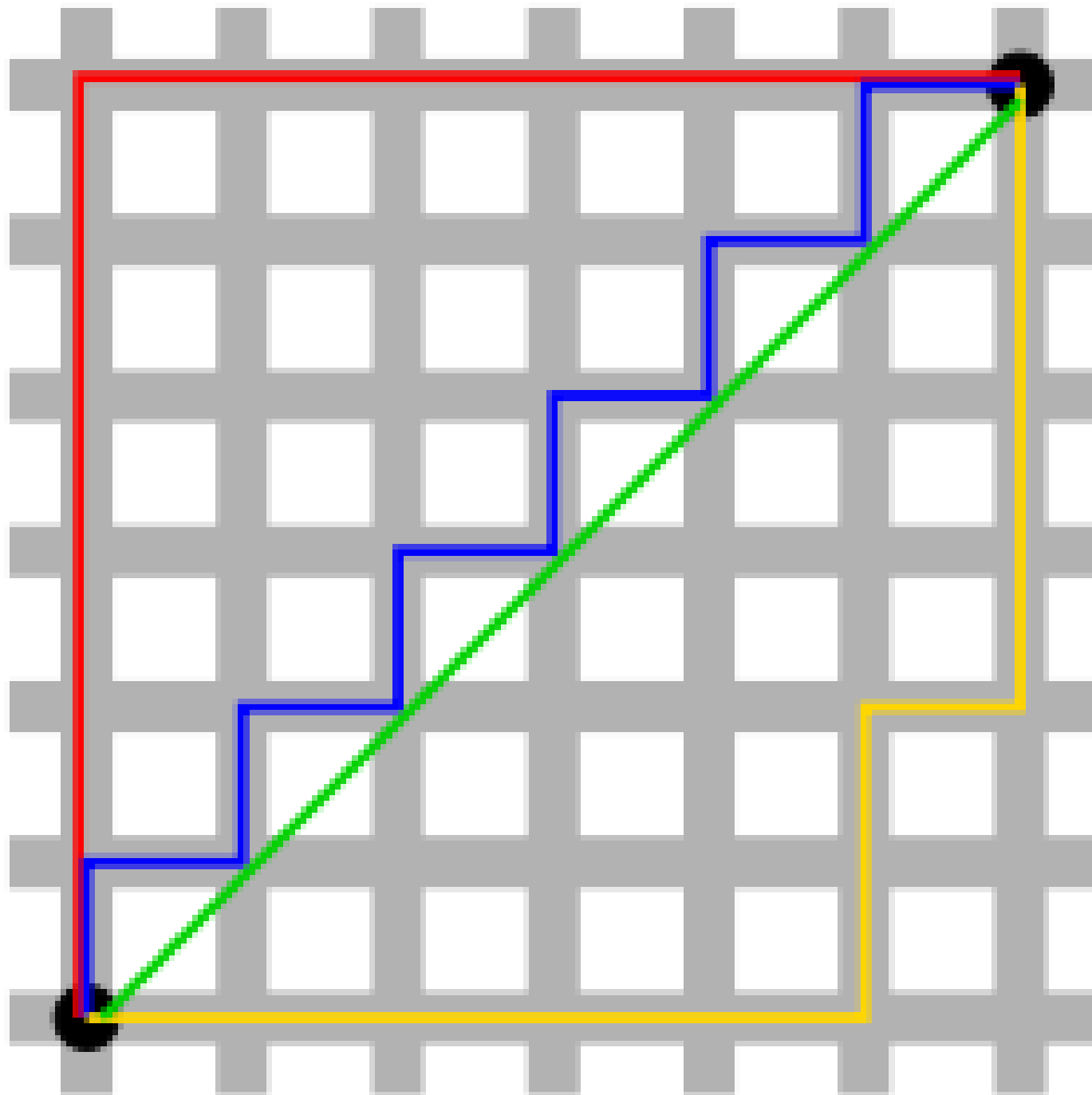
Podobnosti používané u k -NN

- ◆ Často podobnost x a $y \sim \frac{1}{d(x,y)}$ vzdálenost, obvykle $d(x,y) = \|x - y\|$
- 1. $\|x\|_p = \|(x_1, \dots, x_n)\|_p = \sqrt[p]{\sum_{i=1}^n |x_i|^p}$, $p \in [1, +\infty]$ ($p = +\infty$: $\max_i |x_i|$)

$p=1$:

citi-block

vzdálenost





Podobnosti používané u k -NN

◆ Často podobnost x a $y \sim \frac{1}{d(x,y)}$ *vzdálenost*, obvykle $d(x,y) = \|x - y\|$

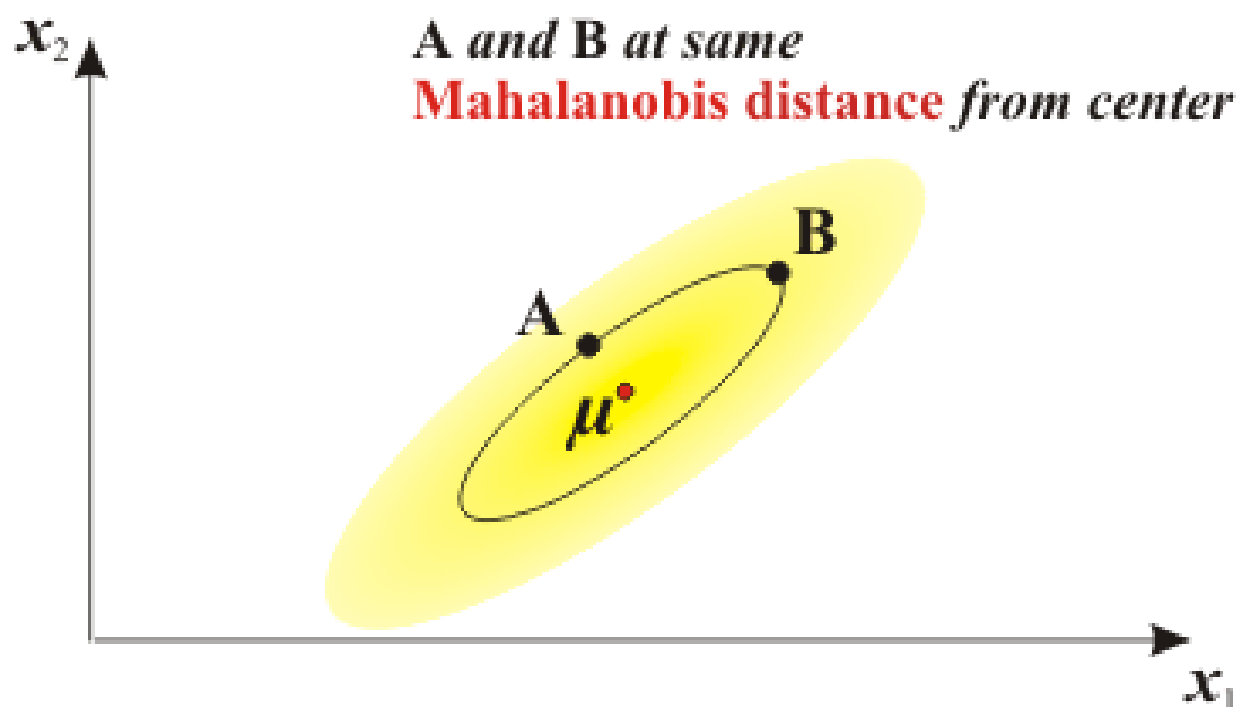
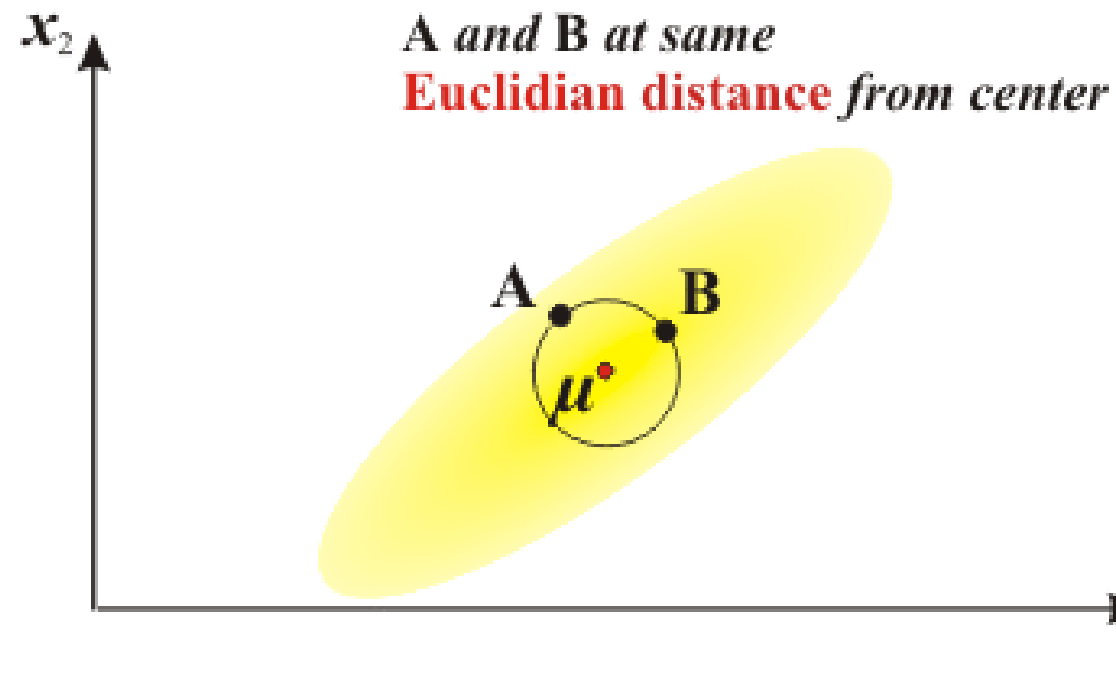
1. $\|x\|_p = \|(x_1, \dots, x_n)\|_p = \sqrt[p]{\sum_{i=1}^n |x_i|^p}$, $p \in [1, +\infty]$ ($p = +\infty$: $\max_i |x_i|$)

2. *Mahalanobisova*: náhodné vektory X, Y , $\text{Var}X = \text{Var}Y = \Sigma$ pozitivně definitní

$$d(x,y) = \left\| \sqrt{\Sigma^{-1}}(X - \mu) - \sqrt{\Sigma^{-1}}(Y - \mu) \right\|_2 = \sqrt{(X - Y)^T \Sigma^{-1} (X - Y)}$$

◆ *Hammingova* $d(x,y) = \#\{i : x_i \neq y_i\}$, *Jaccardova* podobnost $\frac{\#\{i : x_i \neq y_i\}}{\#\{i : \max(x_i, y_i) = 1\}}$

◆ Korelační koeficient x a y , obvykle Pearsonův: $\frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}$



k -NN při kolaborativním filtrování

- ◆ Typicky používané *podobnosti*: spojité atributy – *Pearsonův* korelační koeficient, nominální atributy – *Jaccardova* podobnost
- ◆ Problém: *vychýlenost* (bias) doporučovaného | ostatních uživatelů
 - někteří *uživatelé* hodnotí systematicky přísněji | mírněji
 - některé *produkty* jsou obvykle nadhodnocené | podhodnocené
 - před měřením podobnosti musíme vychýlenost korigovat

k -NN při odhalování malware

- ◆ Používají se *statické vlastnosti* software (= kód),
ne dynamické (= průběh interakcí s OS)
- ◆ Kódování vlastností: *n-gramy bytů* (bigramy, trigramy, ...)
 - vhodně vybraná množina n-gramů: *profil software*
- ◆ K profilu neznámého software hledáme k nejpodobnějších
v databázi známého software (malware + neškodného)

Odhadování pravděpodobné třídy

Odhad pravděpodobnosti náležení bodu x jednotlivým třídám

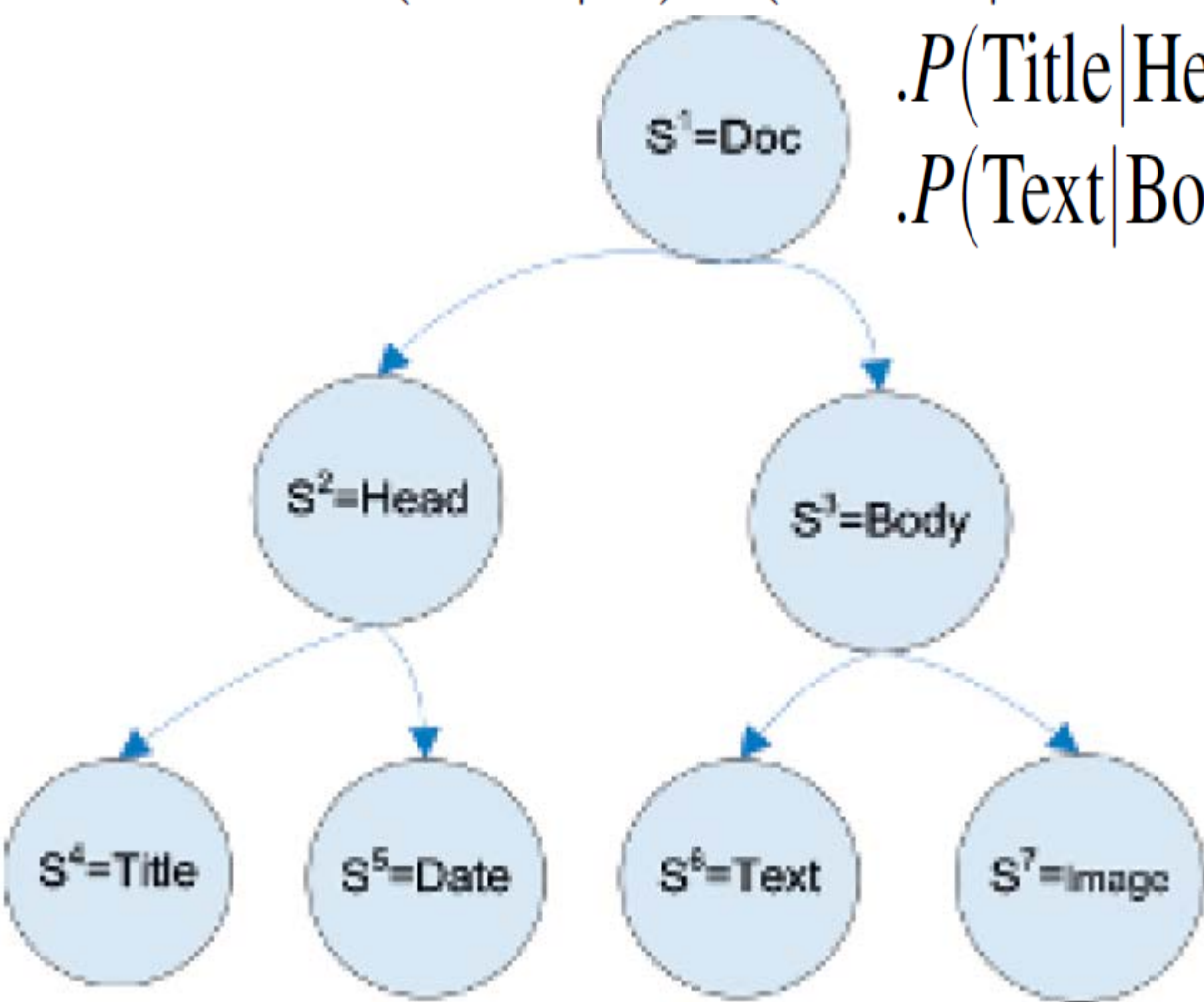
1. Pomocí *Bayesovy věty*: $P(\text{třída } C | \text{atributy } A) = \frac{P(A|C)P(C)}{P(A)}$ \apriorní
 1. naivní klasifikátor: předpokládá nezávislost $P(A_1, \dots, A_n | C) = P(A_1 | C) \cdots P(A_n | C)$
 2. ne naivní: bayesovské sítě (výpočetně náročné)
2. Pomocí *logitového modelu*: jen náležení | nenáležení ($C | \neg C$)

$$P(C | A_1, \dots, A_n) = \frac{1}{1 + e^{-\alpha_1 A_1 \dots - \alpha_n A_n}}, P(\neg C | A_1, \dots, A_n) = \frac{e^{-\alpha_1 A_1 \dots - \alpha_n A_n}}{1 + e^{-\alpha_1 A_1 \dots - \alpha_n A_n}}$$

Bayesovské spamové filtry

- ◆ *Většinou naivní* Bayesův klasifikátor, A_i : slova, metainformace
 - ošálitelné dodatečnými slovy A_i s $P(A_i|\text{ham}) \gg P(A_i|\text{spam})$
 - ⇒ ve výsledku $P(\text{ham}|A_1, \dots, A_n) > P(\text{spam}|A_1, \dots, A_n)$ (*Bayesian poisoning*)
- ◆ Proto alternativně používány jednoduché bayesovské sítě.

$$\begin{aligned}
 P(d) &= P(S^1|\theta)P(S^2|S^1, \theta)P(S^3|S^1, \theta)P(S^4|S^2, \theta) \\
 &\quad .P(S^5|S^2, \theta)P(S^6|S^3, \theta)P(S^7|S^3, \theta) \\
 &= P(\text{Doc}|\theta)P(\text{Head}|\text{Doc}, \theta)P(\text{Body}|\text{Doc}, \theta) \\
 &\quad .P(\text{Title}|\text{Head}, \theta)P(\text{Date}|\text{Head}, \theta) \\
 &\quad .P(\text{Text}|\text{Body}, \theta)P(\text{Image}|\text{Body}, \theta)
 \end{aligned}$$



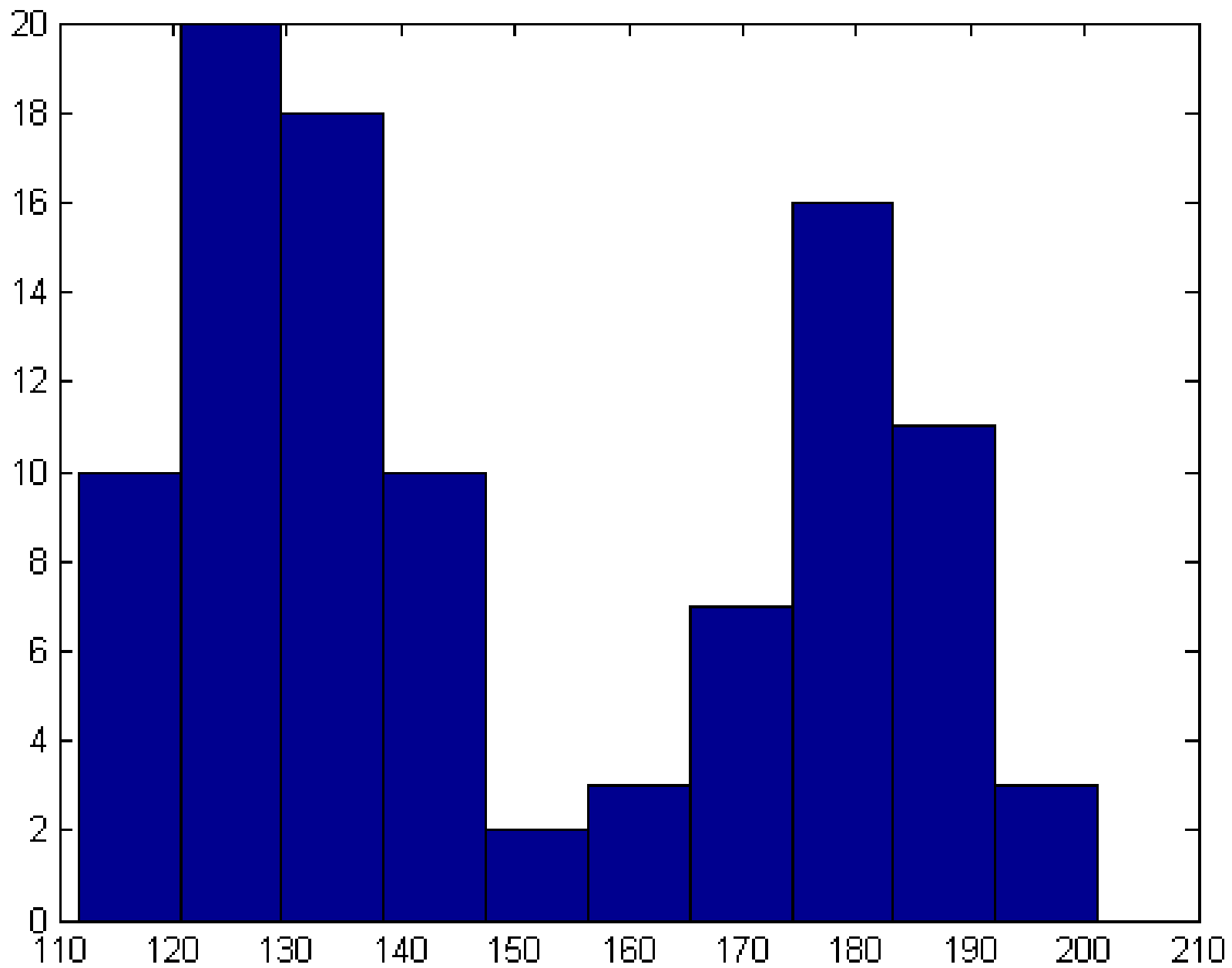
Bayesovské spamové filtry

- ◆ *Většinou naivní* Bayesův klasifikátor, A_i : slova, metainformace
 - ošálitelné dodatečnými slovy A_i s $P(A_i|\text{ham}) \gg P(A_i|\text{spam})$
 - ⇒ ve výsledku $P(\text{ham}|A_1, \dots, A_n) > P(\text{spam}|A_1, \dots, A_n)$ (*Bayesian poisoning*)
- ◆ Proto alternativně používány jednoduché bayesovské sítě.
- ◆ Pro klienty | známější pro *příchozí servery*:

SpamAssassin, DSPAM, Bogofilter, SpamBayes (+ 3. třída: nejasné)

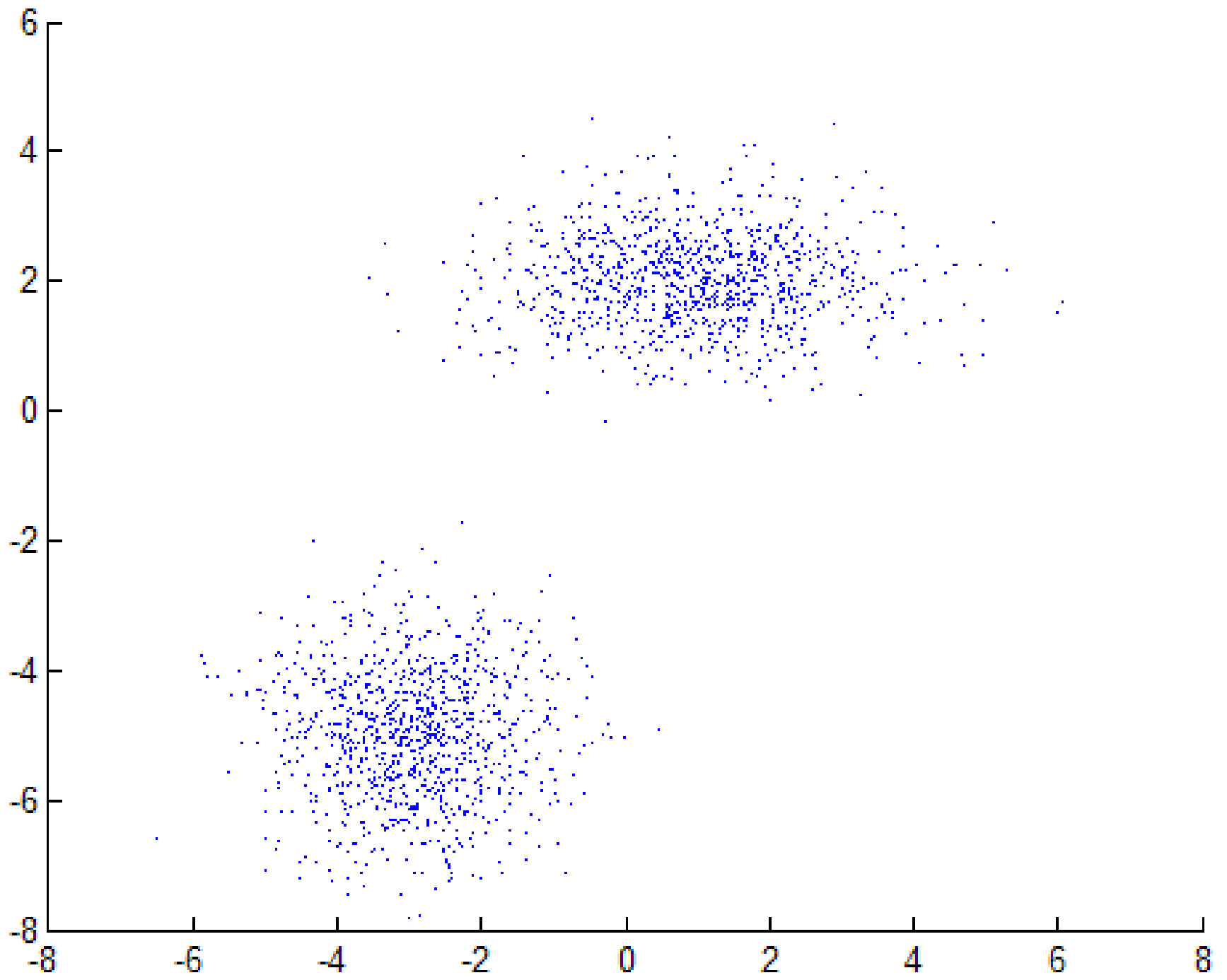
Využití rozdělení pravděpodobnosti

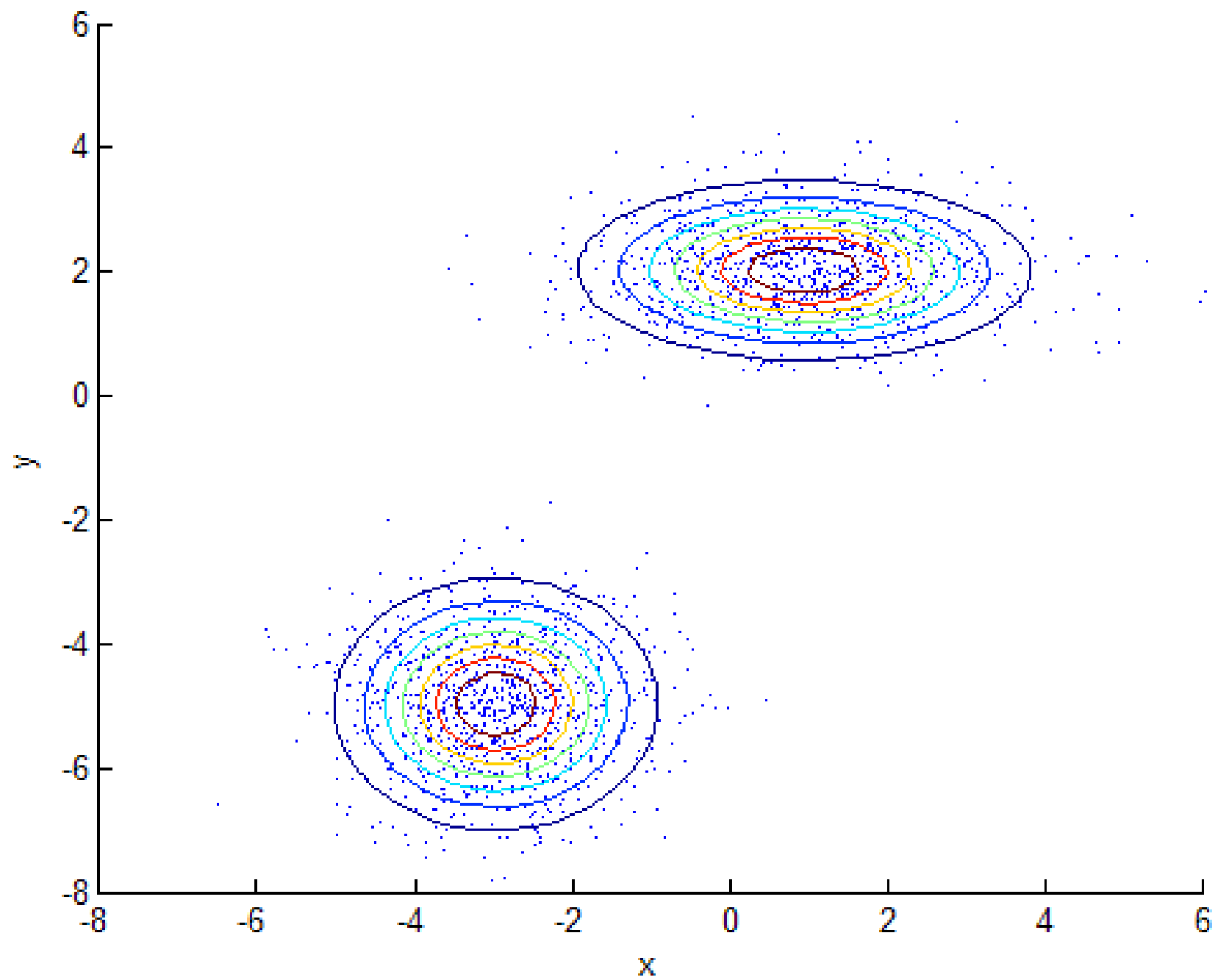
- ◆ *Diskriminační analýza* (DA): $P(C)$ – z rozdělení odpovídajícího C



Využití rozdělení pravděpodobnosti

- ◆ *Diskriminační analýza* (DA): $P(C)$ – z rozdělení odpovídajícího C
 - určeno hustotou $f_C \in \text{rodina } \mathcal{F} = \{f(\cdot|\theta): \theta \in \Theta\}$, Θ – možné parametry
 - originální (Fisherova) DA: \mathcal{F} = vícerozměrná normální rozdělení
- ◆ *Učení* klasifikátoru: výpočet odhadů $\hat{\theta}$ parametrů pro $f_C, C \in \mathcal{C}$
 - $\hat{\theta}$ – maximálně věrohodný odhad: $\hat{\theta} = \arg \max_{\theta \in \Theta} \prod_{x \in \text{učicí data}} f(x|\theta)$
- ◆ Predikce pro nové x : $\arg \max_{C \in \mathcal{C}} f_C(x|\hat{\theta})$





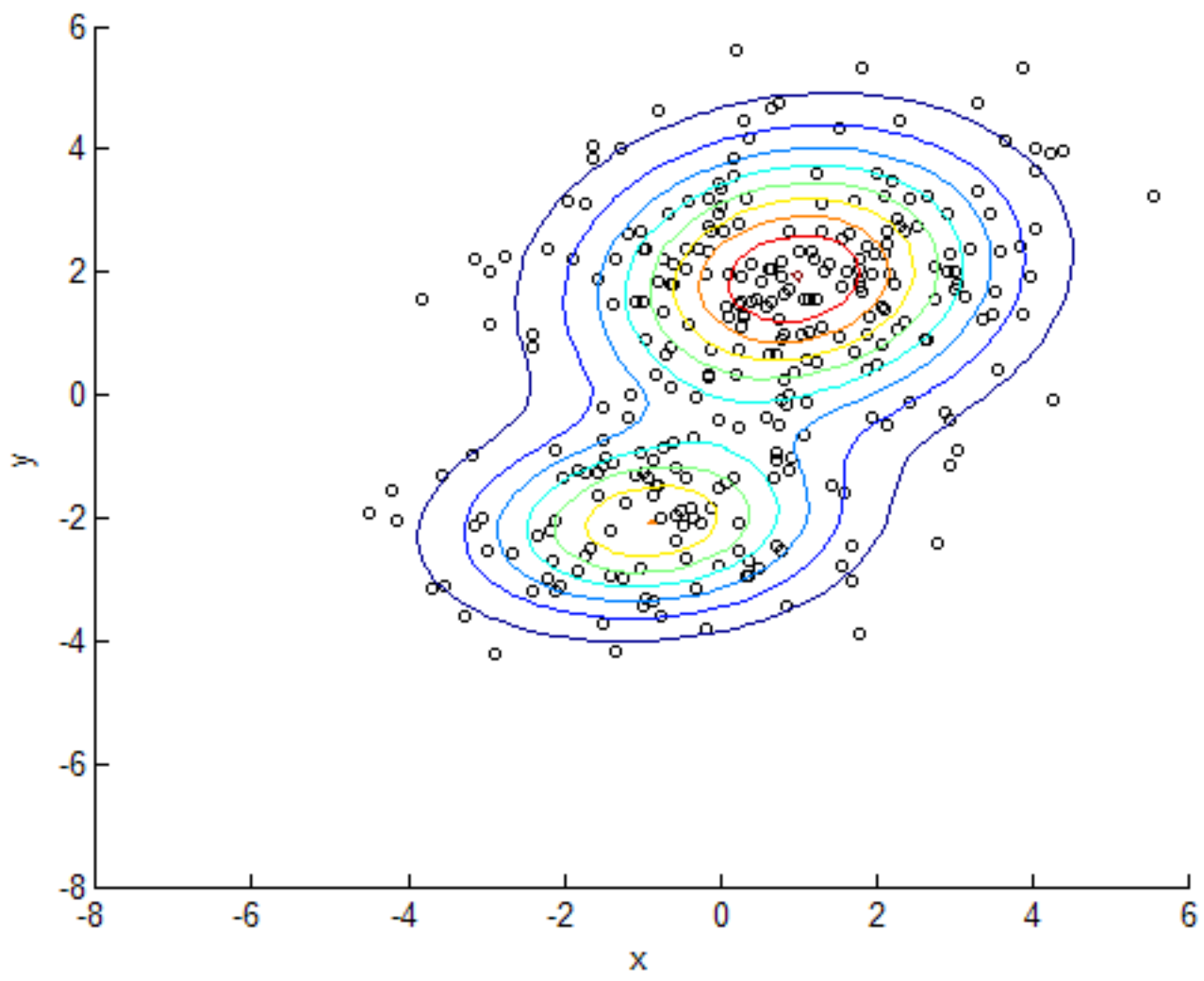
Lineární a kvadratická DA

◆ *Lineární (LDA):* $\mathcal{F} = \left\{ f(x|\beta) = \frac{1}{\sqrt{2\pi|\Sigma|}} e^{-\frac{1}{2}(x-\beta)^\top \Sigma^{-1}(x-\beta)} : \beta \in \mathbb{R}^d \right\}$ dimenze x

- $\Sigma \in \mathbb{R}^{d,d}$: předem daná, stejná pro všechny třídy
- pro $C \neq C'$ je množina $\{x \in \mathbb{R}^d : f_C(x) = f_{C'}(x)\}$ nadrovina (= lineární plocha)

◆ *Kvadratická (QDA):* $\mathcal{F} = \left\{ f(x|\beta) = \frac{1}{\sqrt{2\pi|\Sigma|}} e^{-\frac{1}{2}(x-\beta)^\top \Sigma^{-1}(x-\beta)} : \beta \in \mathbb{R}^d, \Sigma \in \mathbb{R}^{d,d} \right\}$

- Σ – další parametr \Rightarrow parametry celkem $\frac{d(d+1)}{2}$ dimenzí (LDA: d)
- pro $C \neq C'$ je množina $\{x \in \mathbb{R}^d : f_C(x) = f_{C'}(x)\}$ kvadratická plocha



DA při klasifikaci obrázků a videí

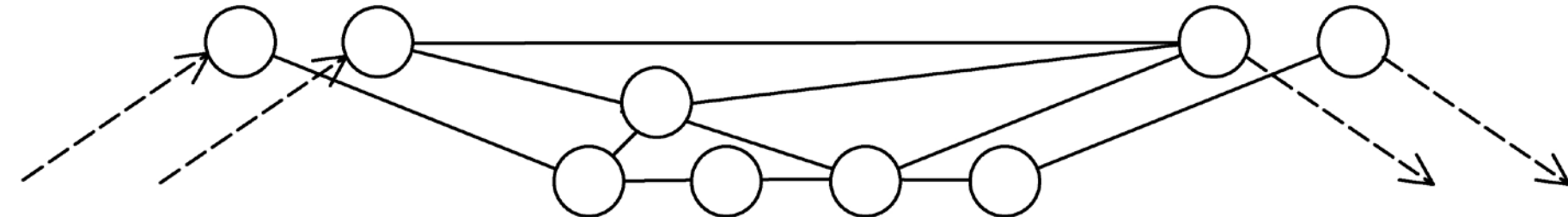
- ◆ Čemu odpovídají jednotlivá pravděpodobnostní rozdělení?
 1. *konceptům* popisujícím zájmy uživatele
 2. *dodaným příkladům* zájmů uživatele
- ◆ *Atributy: statické* (obličeje,...), *dynamické* (pohyby, události)
- ◆ Využití: 1) doporučovací systémy – obsahové filtrování
2) automatická detekce pornostránek

Umělé neuronové sítě

- ◆ Systémy implementující část funkcionality neuronových sítí
- ◆ *Neurony* – vstupní, výstupní, skryté

input neurons

output neurons



incoming signals

hidden neurons

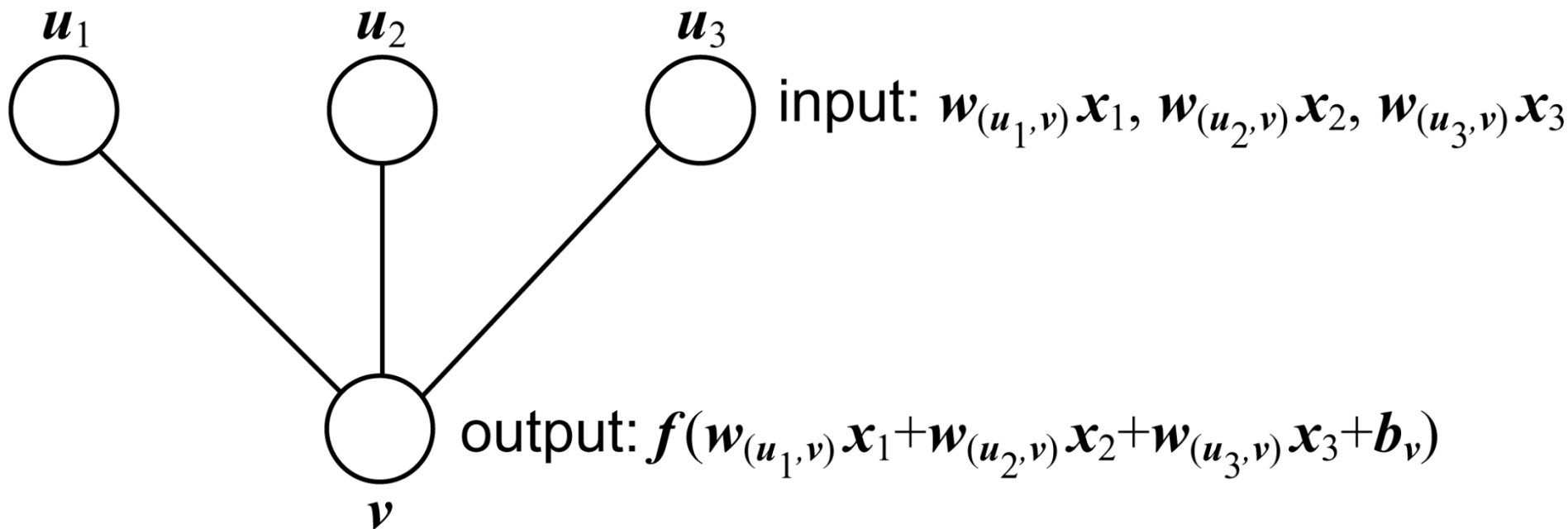
outgoing signals

Umělé neuronové sítě

- ◆ Systémy implementující část funkcionality neuronových sítí
- ◆ *Neurony* – vstupní, výstupní, skryté
 - transformují signály (většinou nelineárně) – *somatické* operace
 - šíření signálů: vstupy → skryté neurony → výstupy
- ◆ *Spoje* – *synaptické* operace: většinou násobení vahou

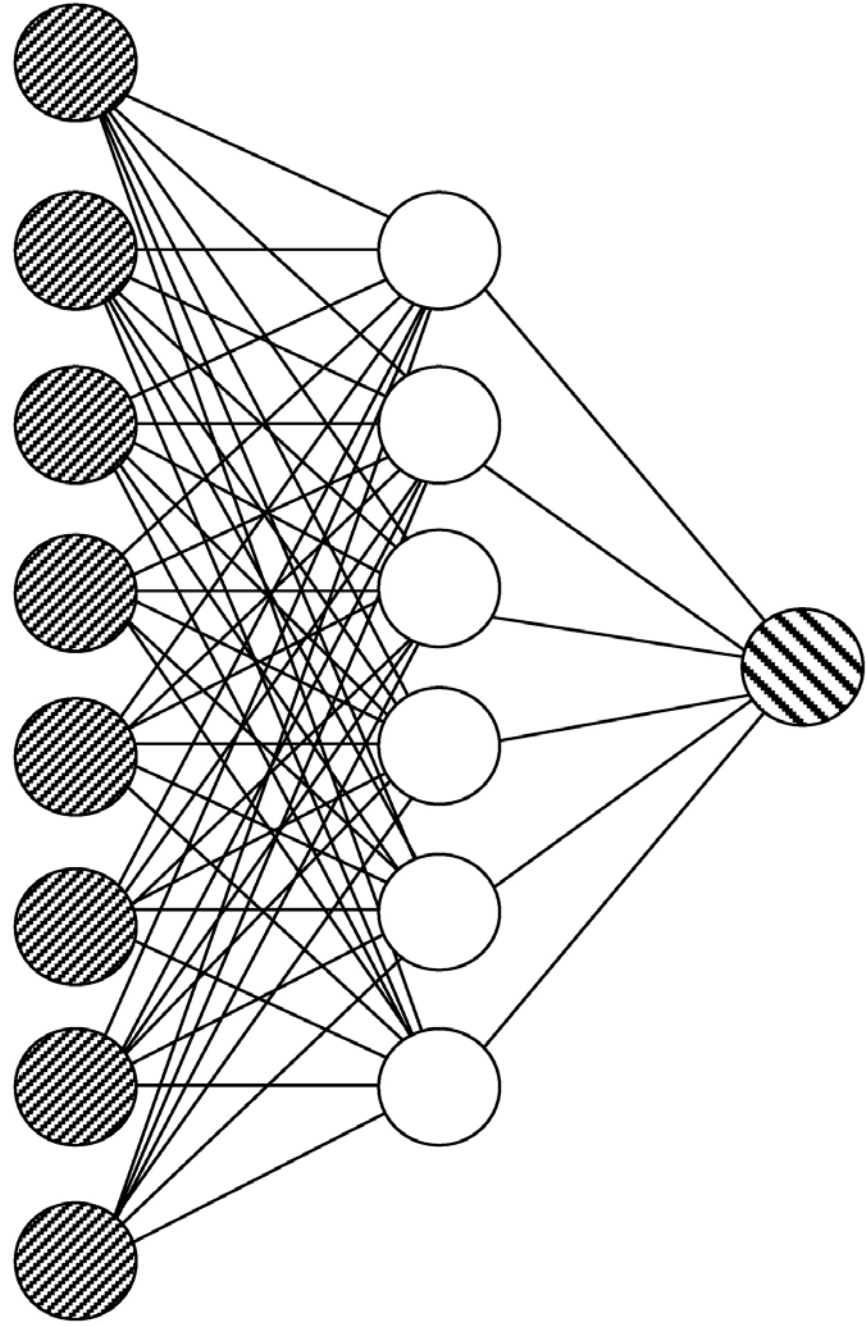
input: \mathbf{x}

output: $\mathbf{W}_{(u,v)} \mathbf{x}$



Umělé neuronové sítě

- ◆ Systémy implementující část funkcionality neuronových sítí
- ◆ *Neurony* – vstupní, výstupní, skryté
 - transformují signály (většinou nelineárně) – *somatické* operace
 - šíření signálů: vstupy → skryté neurony → výstupy
- ◆ *Spoje* – *synaptické* operace: většinou násobení vahou
- ◆ *Topologie* – propojení neuronů spoji: nejčastěji *vrstevnatá*



input layer

hidden layer

output layer

Perceptrony a jejich zřetězení

- ◆ *Perceptron* – bez skrytých neuronů, somatická operace:

skoková aktivační funkce Θ : $\Theta(\mathbb{R}_-) = 0, \Theta(\mathbb{R}_{+0}) = 1 \Rightarrow f_C(x) = \Theta(w^T x + b)$

- třída \subset průnik poloprostorů \Rightarrow vyžaduje lineární separabilitu

- ◆ *Vícevrstvý perceptron (MLP)* – zřetězení \Rightarrow skryté vrstvy

- somatická operace: sigmoidní aktivační funkce ($\sigma(x) = \frac{1}{1+e^{-x}}, = \frac{e^x - e^{-x}}{e^x + e^{-x}}, \dots$)

- nelineární klasifikátory, při 1 skryté $f_C = \sigma(\sum_{h \in H} w_{hc} \sigma(w_{.h}^T x + b_h) + b_c)$

MLP při filtraci spamu

- ◆ Výhoda: vysoká *přesnost*, zejména týmového použití
- ◆ Nevýhoda: *pomalé učení* (+ doučování nových spamů)
- ◆ ⇒ *Strategie* pro nasazení ve spamových filtrech:
 - na straně *klienta* (⇐ doučování méně časté)
 - doplňovány klasifikátory s rychlým doučováním (k-NN, ...)
- ◆ Další použití: filtrace spamu přenášeného *obrázky*

MLP v doporučovacích systémech

- ◆ Při kolaborativním filtrování častěji než při obsahovém
- ◆ Výhoda: *přesnost*, zase hlavně v týmu
- ◆ Nevýhoda: *obtížná vysvětlitelnost* (> než pomalé učení)

výsledku: proč je konkrétní objekt doporučován

- ◆ Alternativní použití: *predikce nesdělených* sociálních *atributů* (věk, příjmová skupina, ...) indikujících, co doporučit

MLP při odhalování malware

- ◆ Odhalují podle statických vlastností, tj. kódu
 - vstupy: n-gramy bytů | operací, metadata binárek
 - kombinování s klasifikátory odhalujícími podle chování
- ◆ Problém: špatné učení třídy zastoupené vzácně
 - protiopatření: obohacení učicích dat zástupci malware
- ◆ Příklad komerčního použití MLP: antivirák IBM