

Internet a klasifikační metody, 6. přednáška

Tým zvládne víc než jedinec

volitelný předmět pro magisterské studium

Martin Holeňa



O čem to bude?

- ◆ Spojování více klasifikátorů do týmu
 - zahrnutí různé důvěry různým členům týmu
 - týmy klasifikátorů různých a stejných typů
 - hlavní metody spojování: bagging, boosting, stacking
- ◆ Hierarchické uspořádání týmů klasifikátorů
- ◆ Týmy klasifikačních stromů – klasifikační náhodné lesy

Spojování klasifikátorů do týmu

- ◆ *Motivace*: máme-li tým klasifikátorů, nepotřebujeme aby každý klasifikoval správně jakýkoliv vstup
 - *možnost specializace* na část prostoru příznaků X
 - + nutnost odhadnout důvěry jednotlivým členům týmu
- ◆ Pro r -členný tým: $F: X \rightarrow \{C_1, \dots, C_m\}, F(x) = A(F_1(x), \dots, F_r(x), \tau_1(x), \dots, \tau_r(x))$
 - A – *agregační funkce*, τ_i – *důvěra* (trust) klasifikátoru F_i , $\tau_i: X \rightarrow [0,1]$

Jak měřit důvěru klasifikátoru

1. *Globální* důvěra klasifikátoru F – nezávisí na x , $(\forall x \in X) \tau(x) = \tau_F \in [0,1]$

- příklad $\tau_F = \frac{1}{|V|} \sum_{v \in V} \|F(v) - C_v\|$, kde $v \in C_v, V$ – daná validační množina

2. *Lokální* důvěra – závisí na x , nejznámější příklad:

Euclidean Local Accuracy – používá eukleidovskou vzdálenost d_E

a k nejbližších sousedů z V , $\tau^{ELA}(x) = \frac{1}{k} \sum_{i=1}^k \|F(v^{(x,i)}) - C_{v^{(x,i)}}\|$

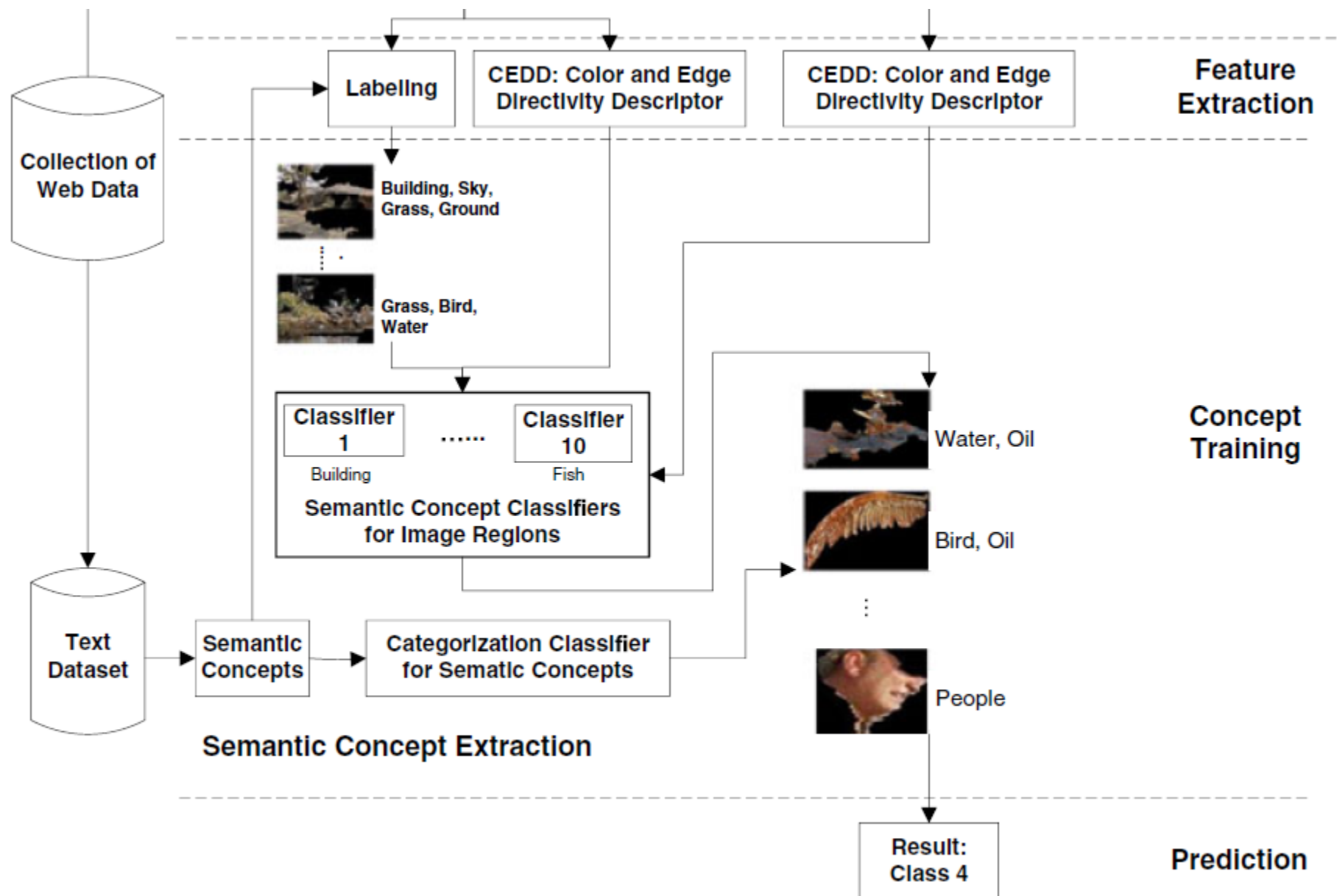
- $v^{(x,1)}, \dots, v^{(x,|V|)}$ – uspořádání prvků V dle vzdálenosti $d_E(v, x)$ od x

Týmy a soubory klasifikátorů

- ◆ *Jak* vnést mezi členy týmu *nepodobnost*?
 - proč potřeba? jinak tým = kterýkoliv člen
- 1. Klasifikátory *různých druhů* (k -NN, Bayesův, SVM, ...)
- 2. Klasifikátory 1 druhu s různými hodnotami parametrů
 - k -NN: k , SVM: c, q v $\kappa(x, y) = (x^T y + c)^q$, MLP: počet neuronů
 - takový tým nazýváme soubor (*ensemble*) *klasifikátorů*

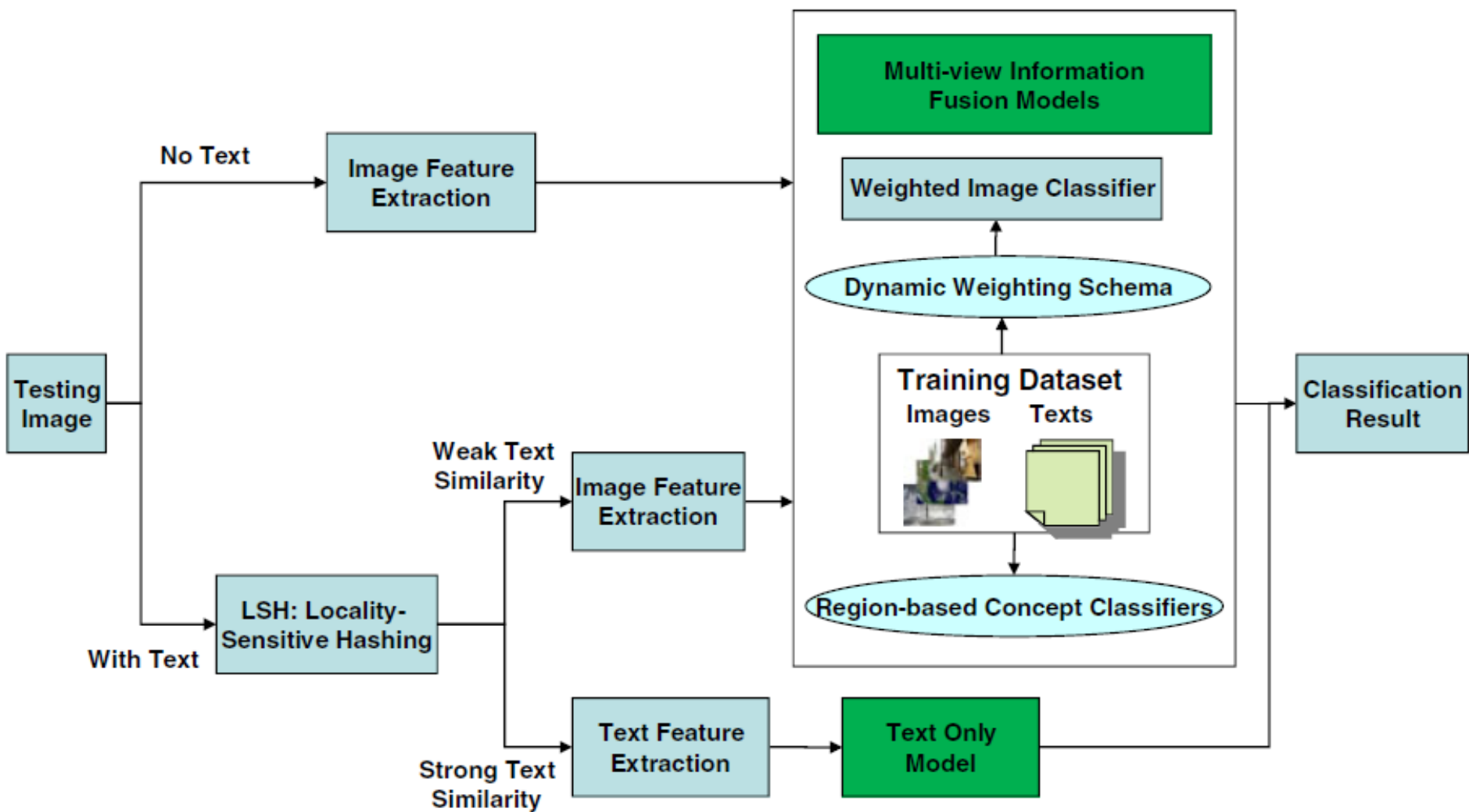
Týmová klasifikace multimedií

- ◆ K plnému využití informací klasifikujeme každé použité medium: obrázky, text, videa, zvuky, ...



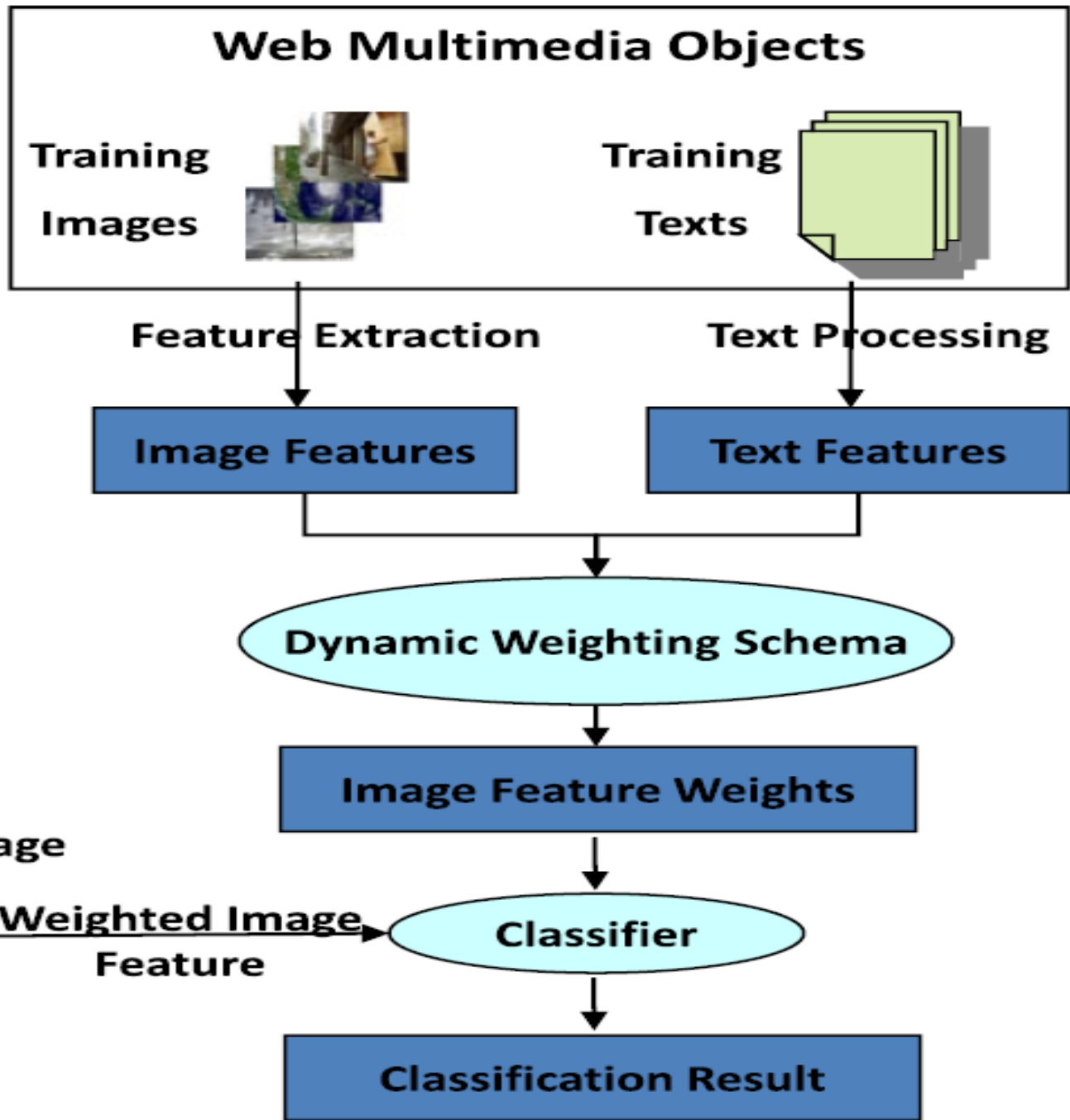
Týmová klasifikace multimedií

- ◆ K plnému využití informací klasifikujeme každé použité medium: obrázky, text, videa, zvuky, ...
- ◆ Různá media nejlépe klasifikována různými typy klasifikátorů → potřebujeme tým *klasifikátorů různých typů*



Týmová klasifikace multimedií

- ◆ K plnému využití informací klasifikujeme každé použité medium: obrázky, text, videa, zvuky, ...
- ◆ Různá media nejlépe klasifikována různými typy klasifikátorů → potřebujeme tým *klasifikátorů různých typů*
- ◆ AgregáčnÍ funkce: klasifikátory odpovídající jednotlivým mediím *dynamicky vážíme důvěrou* (dynamicky – lokální důvěrou)



Web Multimedia Objects

Training Images



Training Texts



Feature Extraction

Text Processing

Image Features

Text Features

Dynamic Weighting Schema

Image Feature Weights

Testing Image



Weighted Image Feature

Classifier

Classification Result

Jak tým klasifikátorů vytvořit?

- ◆ *Různodruhové týmy*: nemnoho členů (nejčastěji jednotky)
 - každý člen naučen nezávisle, specifickou metodou
- ◆ *Soubory*: desítky – tisíce členů, principiálně vytvářeny učením různými podmnožinami dat – bagging, boosting

Bagging

- ◆ **Bootstrap aggregating**: agregované $F_k, k = 1, \dots, r$ trénovány na $\{x_1^{(k)}, \dots, x_n^{(k)}\}$ získaných bootstrapem = *resamplováním s vracením*
 - bootstrap z n dat: počet různých $\rightarrow n \lim_{n \rightarrow \infty} P((\exists j) x_j^{(k)} = x_i) = n \left(1 - \frac{1}{e}\right)$
- ◆ **Validace** pro F_k : daty $x \notin \{x_1^{(k)}, \dots, x_n^{(k)}\}$ („out of bag“)
 - baggingová alternativa křížové validace: $q \geq 2$ částí (q -fold), každá validuje klasifikátor naučený $q - 1$ zbývajícími, + zprůměrování

Boosting

◆ Členové učení všemi učitími daty U , *vliv* $u \in U$

vážen neúspěšností klasifikace u některými jinými členy

- algoritmy: AdaBoost (nejužívanější), LPBoost, LogitBoost, BrownBoost

◆ *AdaBoost* (Adaptive Boosting): klasifikátory $F_i: X \rightarrow \{1, -1\}, i = 1, \dots, r$

- váhy: $(\forall u \in U) w_1(u) = \frac{1}{|U|}, w_{i+1}(u) = w_i(u) e^{\tau_i(2\|F_i(u) \neq C_u\| - 1)}$

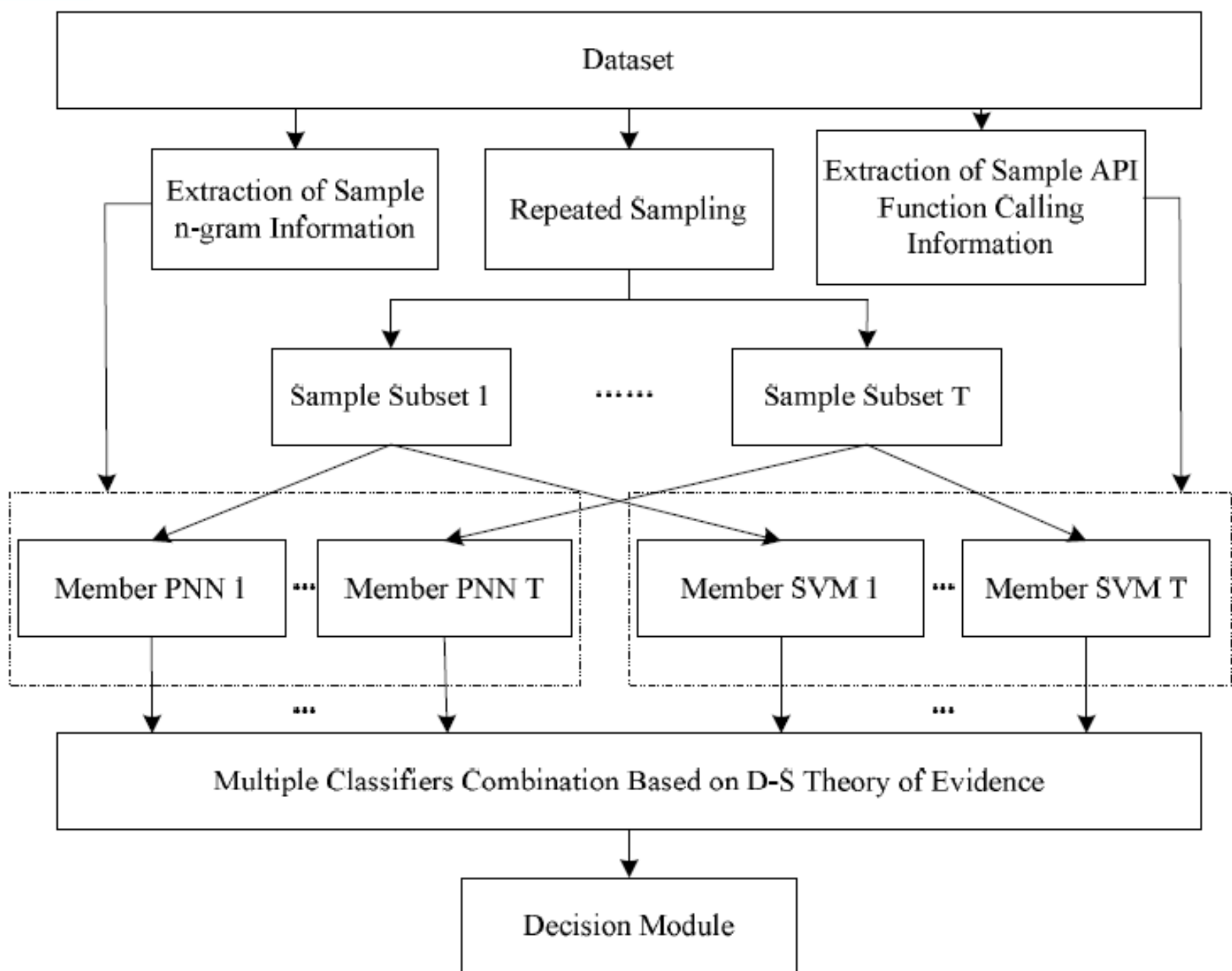
- agregace: $A(F_1, \dots, F_r, \tau_1, \dots, \tau_r) = \text{sign}(\sum_{i=1}^r \tau_i F_i), \tau_i = \frac{1}{2} \ln \frac{\sum_{u \in U} w_i(u) \|F_i(u) = C_u\|}{\sum_{u \in U} w_i(u) \|F_i(u) \neq C_u\|}$

Jak tým klasifikátorů vytvořit?

- ◆ *Různodruhové týmy*: nemnoho členů (nejčastěji jednotky)
 - každý člen naučen nezávisle, specifickou metodou
- ◆ *Soubory*: desítky – tisíce členů, principiálně vytvářeny učením různými podmnožinami dat – bagging, boosting
- ◆ *Dodatečná metoda* vytváření týmů: použít jako člen tým (např. soubor) → *hierarchické* týmy

Soubory klasifikátorů a malware

- ◆ Detekce malware na základě statických vlastností
 - klasifikace *n-gramů* souborem specifických *neuronových sítí*
- ◆ Detekce malware na základě dynamických vlastností
 - klasifikace *volání API* souborem *SVM* klasifikátorů
- ◆ Různí členové souborů trénování různými podmnožinami dat



Soubory klasifikátorů a malware

- ◆ Detekce malware na základě statických vlastností
 - klasifikace *n-gramů* souborem specifických *neuronových sítí*
- ◆ Detekce malware na základě dynamických vlastností
 - klasifikace *volání API* souborem *SVM* klasifikátorů
- ◆ Různí členové souborů trénování různými podmnožinami dat
- ◆ Princip agregace: teorie evidence (Dempster-Shafer)

Náhodné lesy

- ◆ Náhodný les: *soubor stromů*, trénován každý zvlášť, při klasifikaci jsou *predikce sčítány*

Training

$\{v\}$



node 0

\mathcal{S}_1

1

2

\mathcal{S}_1^L

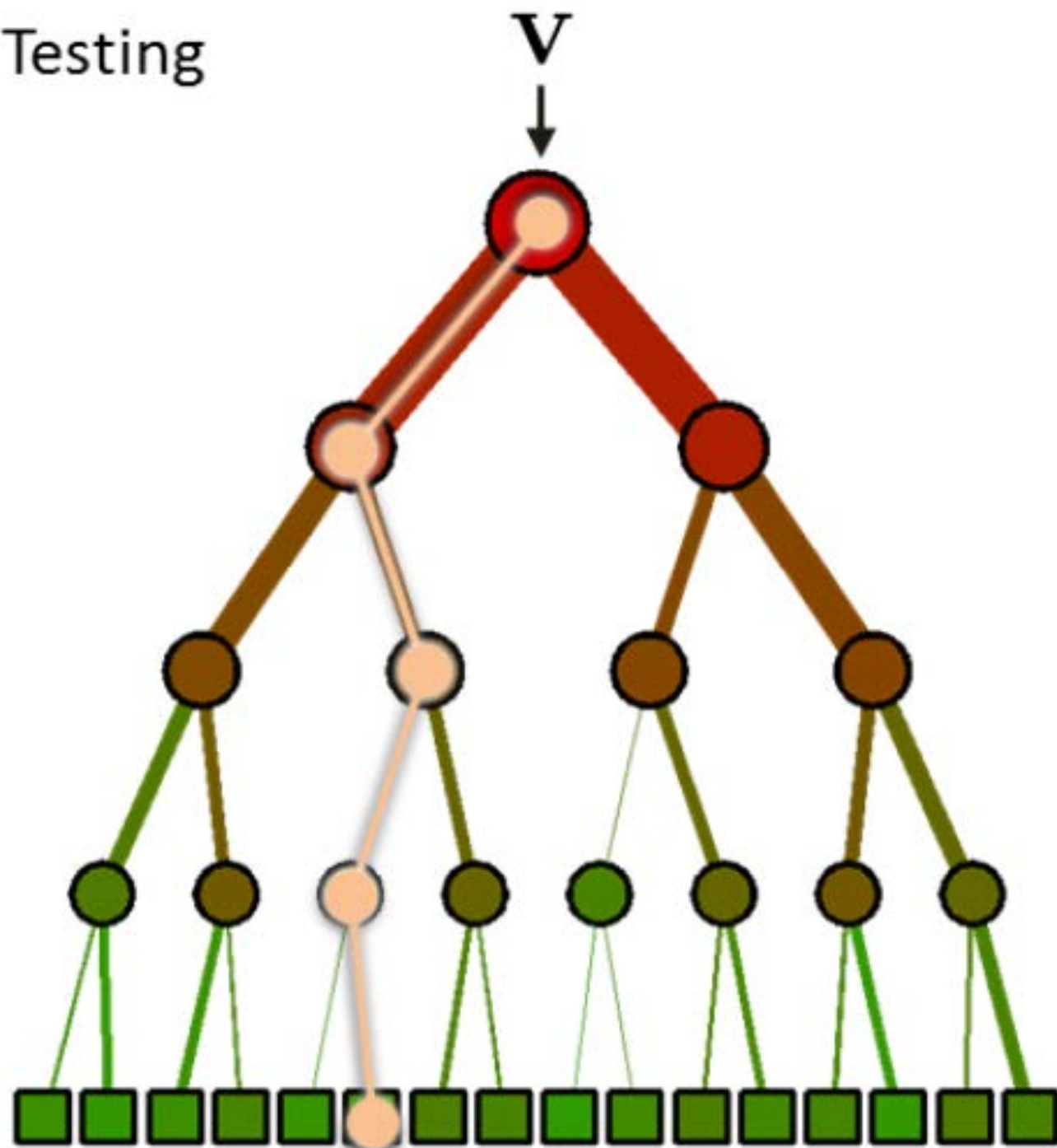
\mathcal{S}_1^R

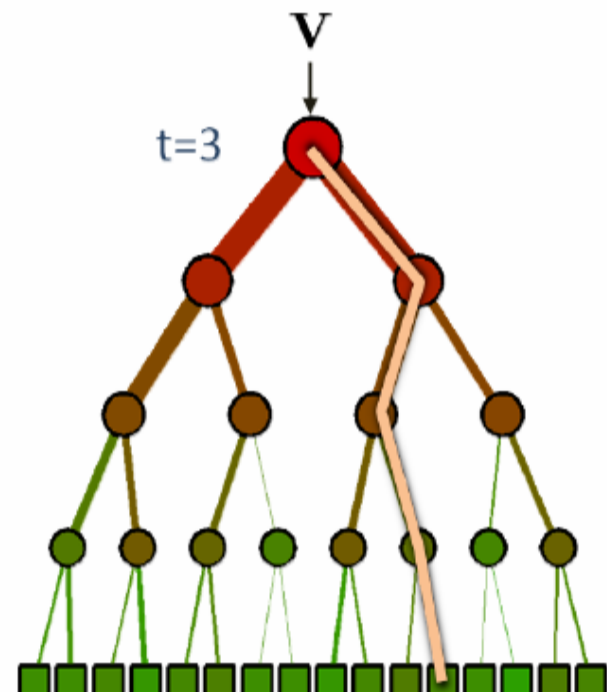
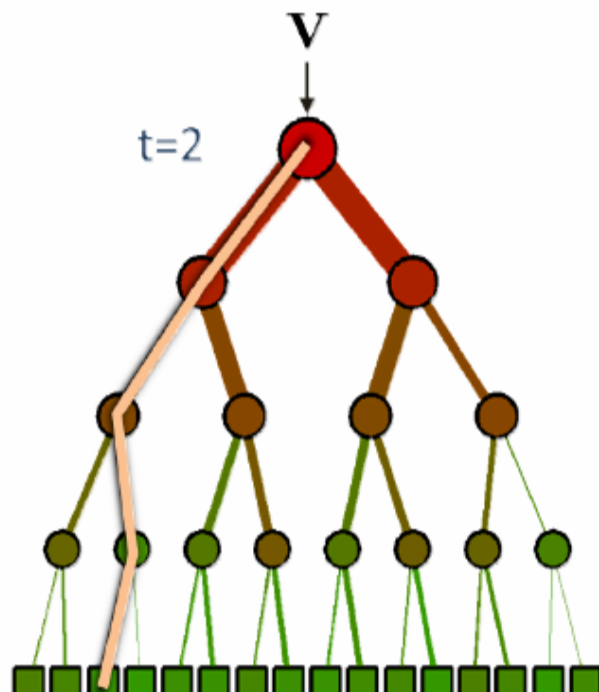
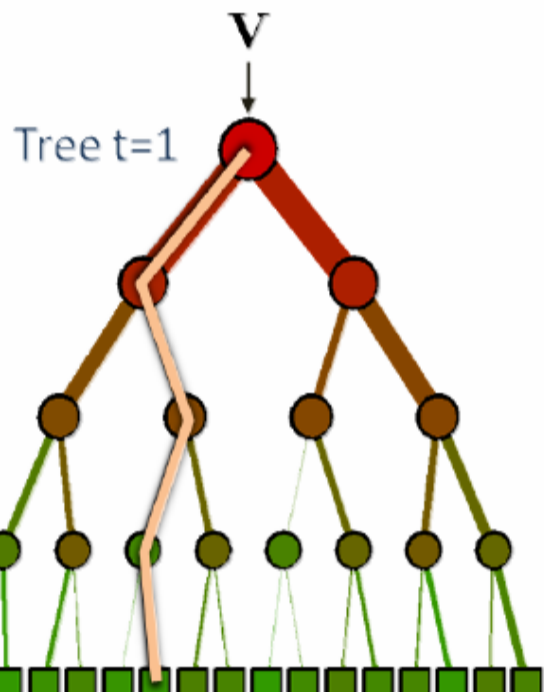
3

4



Testing





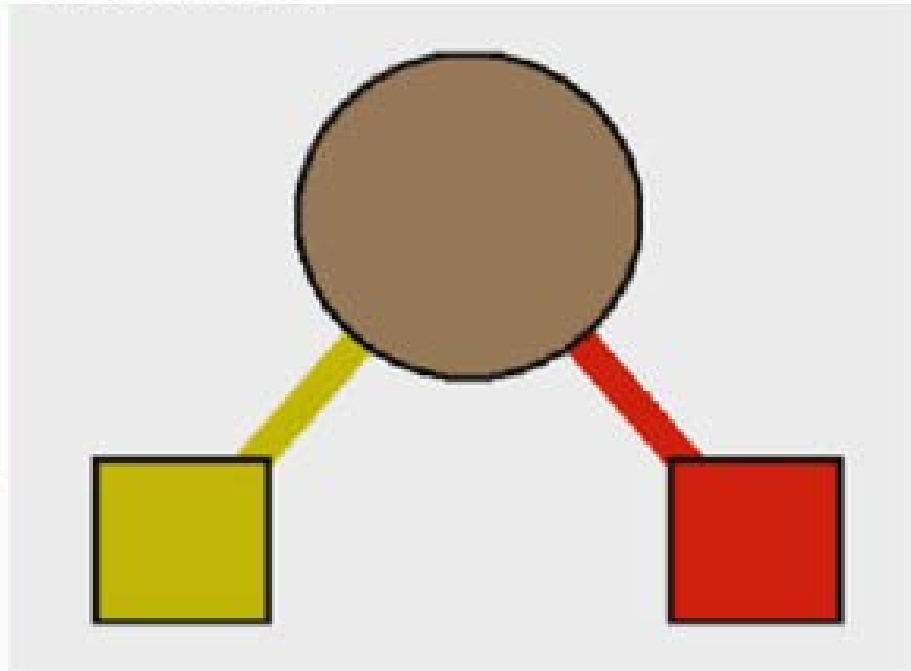
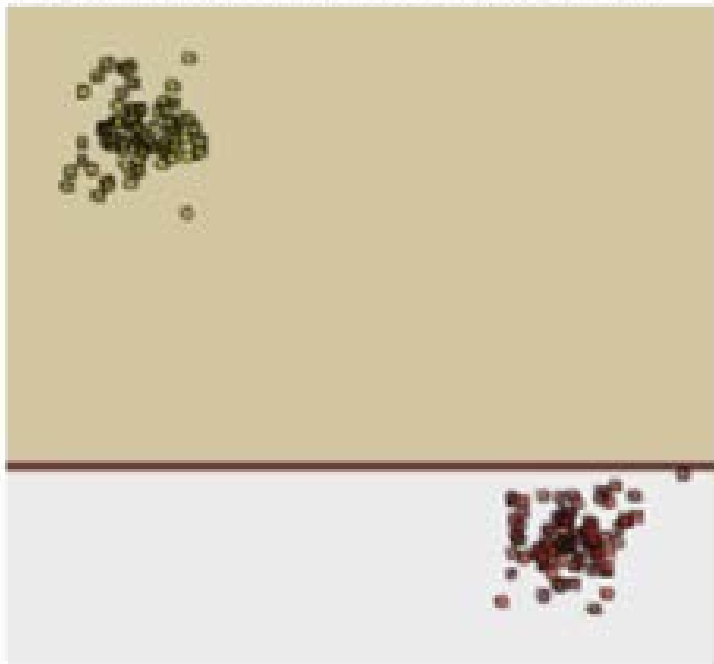
Náhodné lesy

- ◆ Náhodný les: *soubor stromů*, trénován každý zvlášť, při klasifikaci jsou *predikce sčítány*
 - hranice: stromy – ostré, les – postupně přecházející

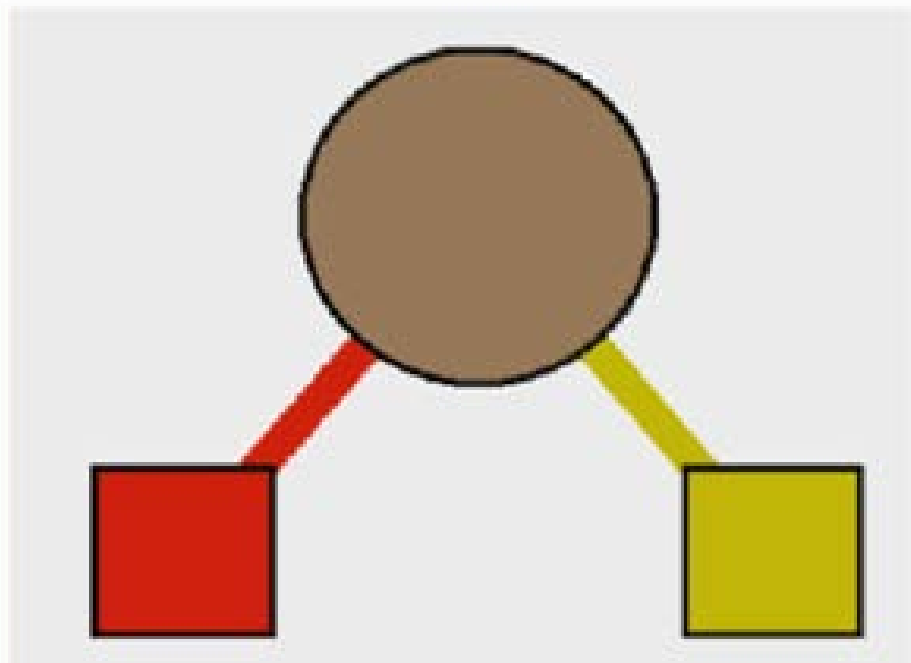
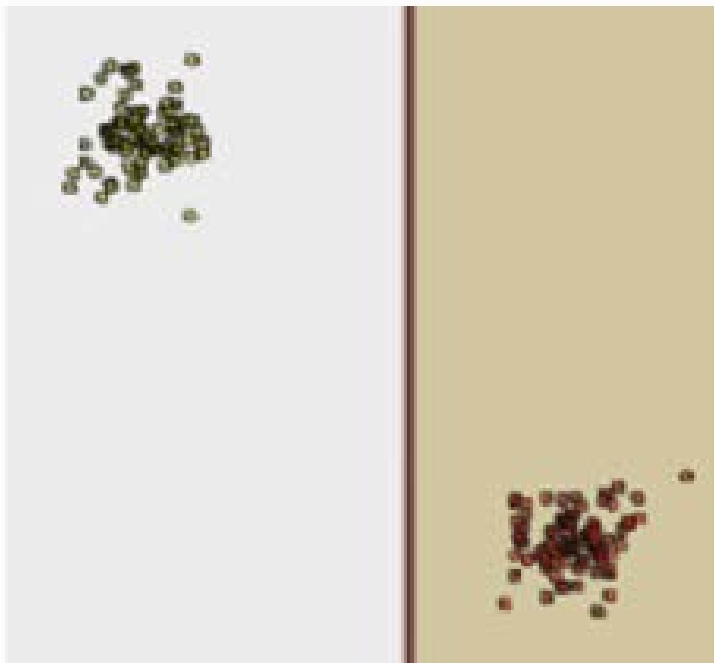
Learnt decision boundaries

Learnt trees

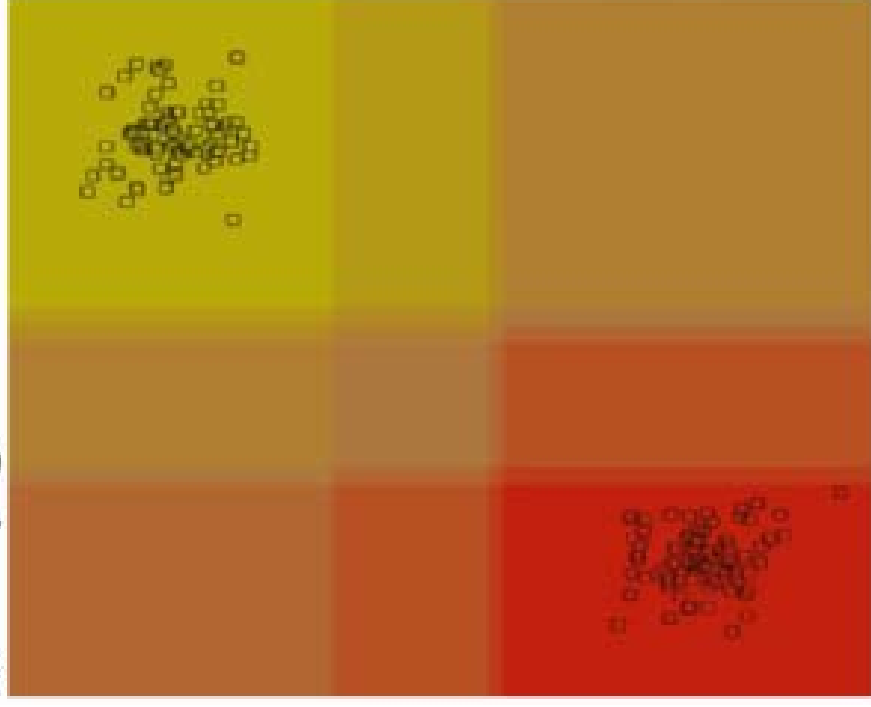
Tree 43



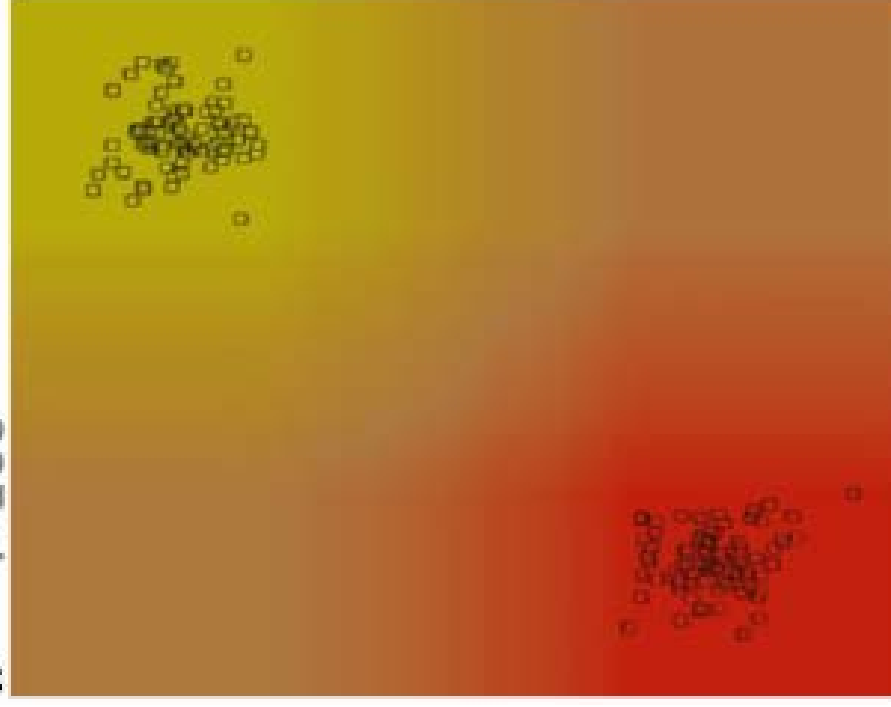
Tree 121



C_2 $T=8$



C_3 $T=200$

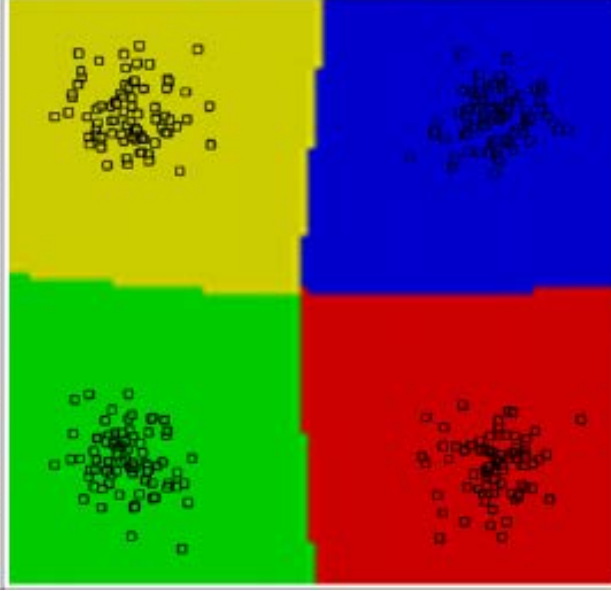


Náhodné lesy

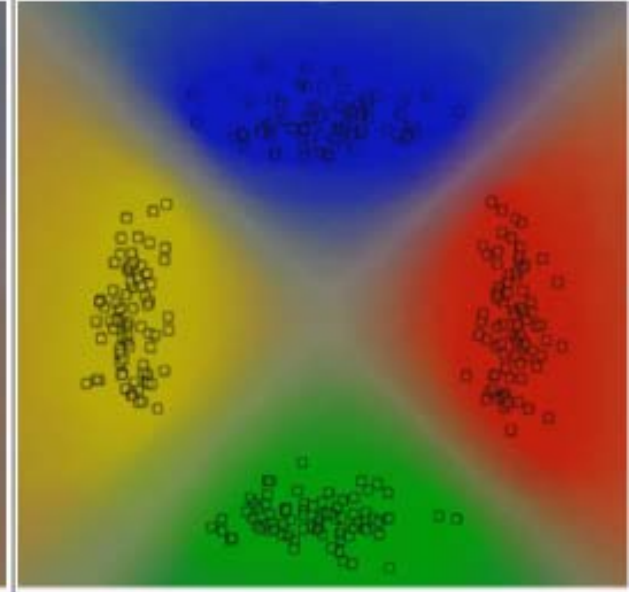
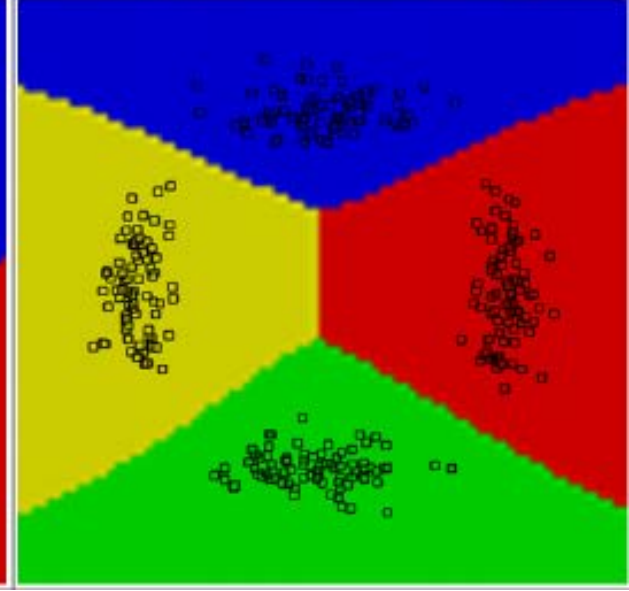
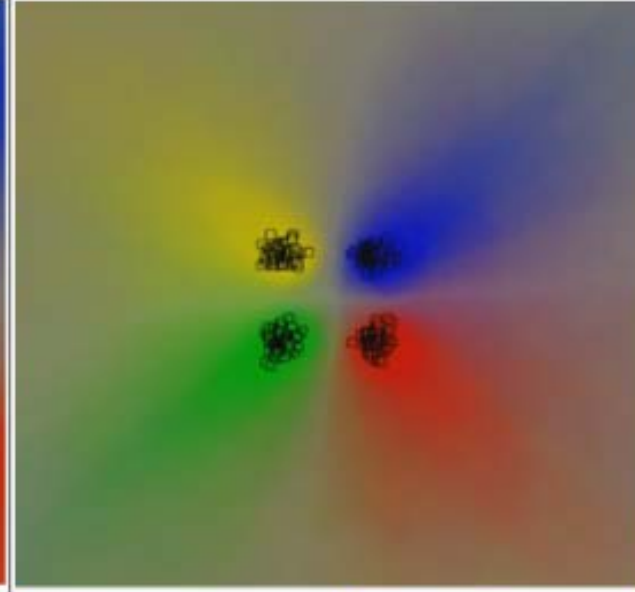
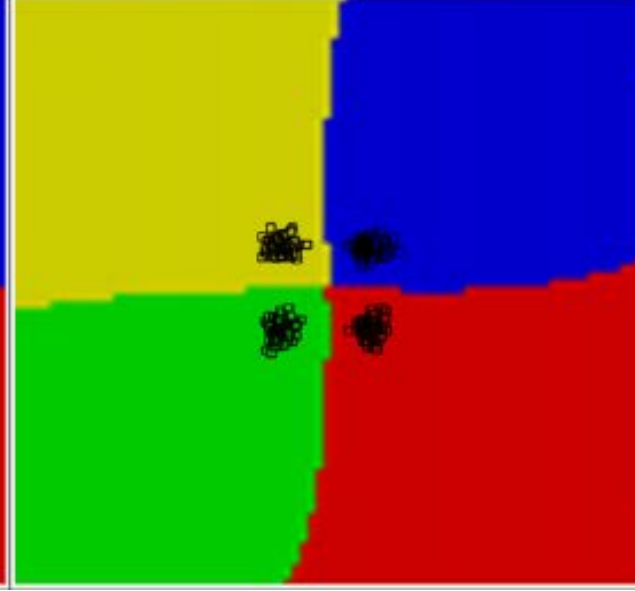
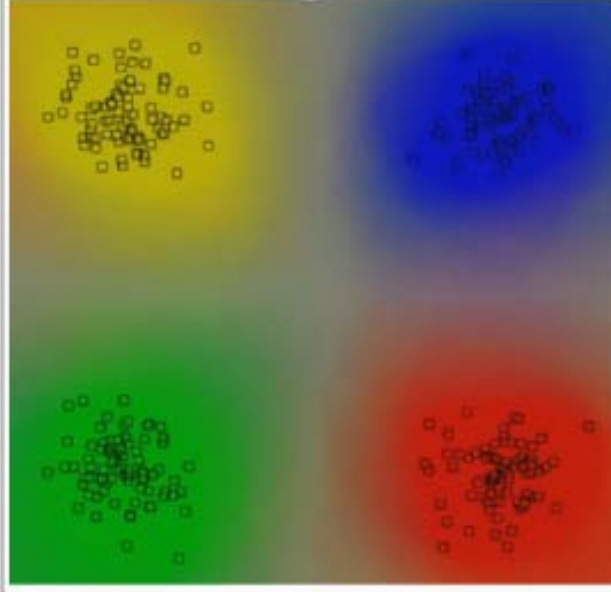
- ◆ Náhodný les: *soubor stromů*, trénován každý zvlášť, při klasifikaci jsou *predikce sčítány*
 - hranice: stromy – ostré, les – postupně přecházející

1. výhoda: pokud třídy ve skutečnosti fuzzy, popisují je přesněji než ostatní klasifikátory

Support vector machine (1-v-all)

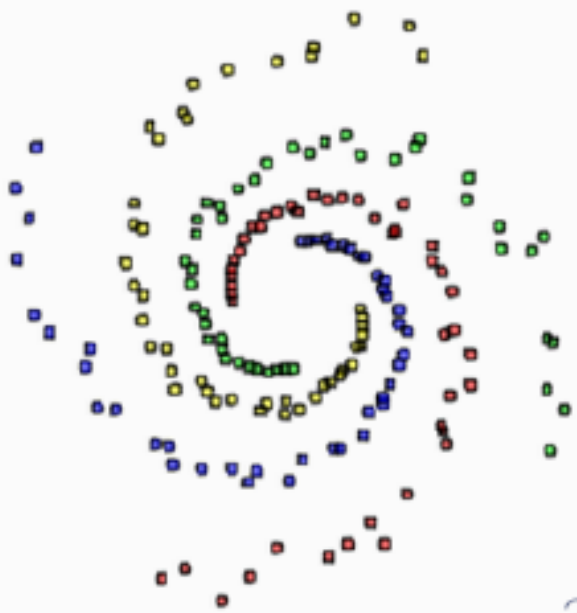
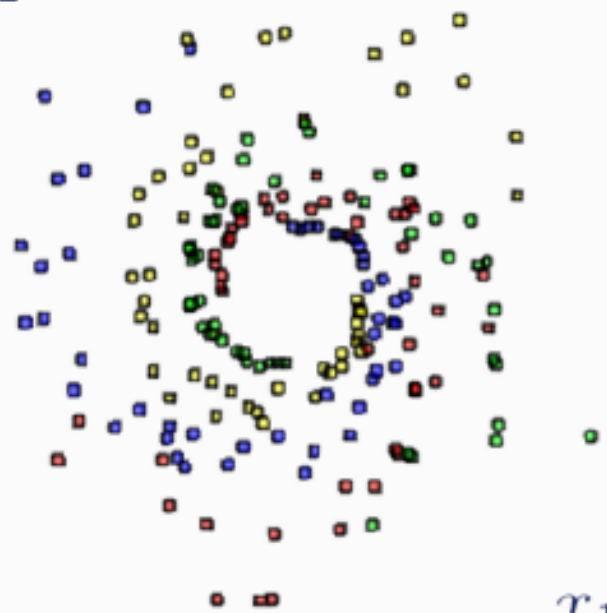
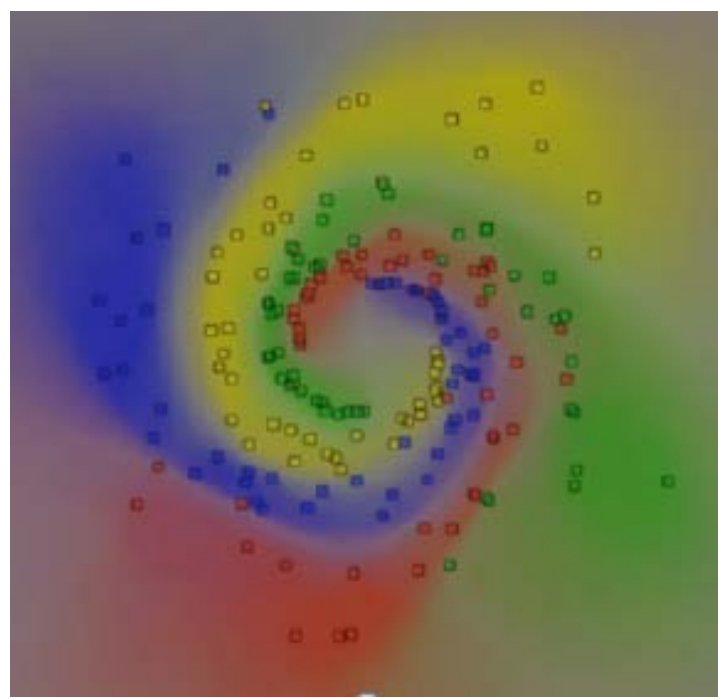
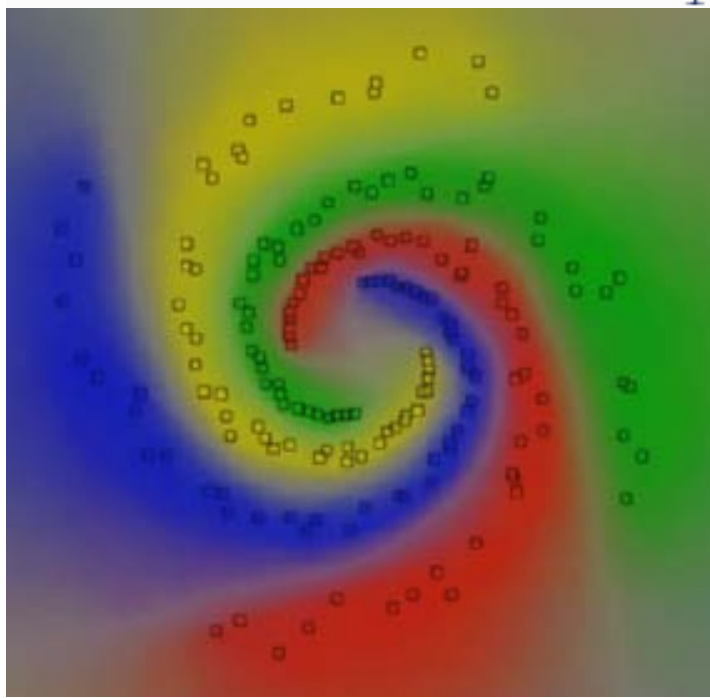


Classification forest



Náhodné lesy

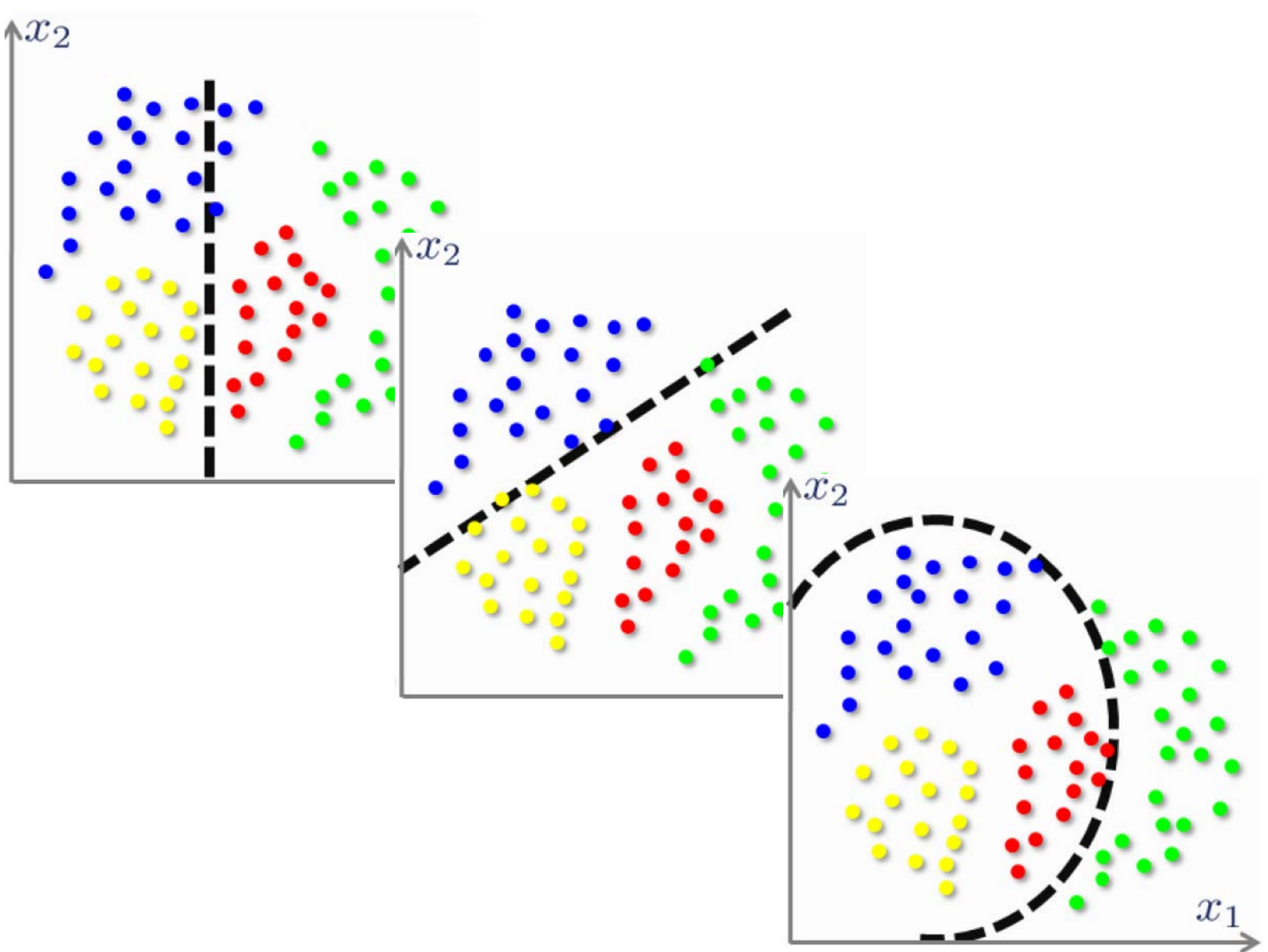
- ◆ Náhodný les: *soubor stromů*, trénován každý zvlášť, při klasifikaci jsou *predikce sčítány*
 - hranice: stromy – ostré, les – postupně přecházející
- 1. *výhoda*: pokud třídy ve skutečnosti fuzzy, popisují je přesněji než ostatní klasifikátory
- 2. *výhoda*: sčítané predikce méně ovlivněny šumem

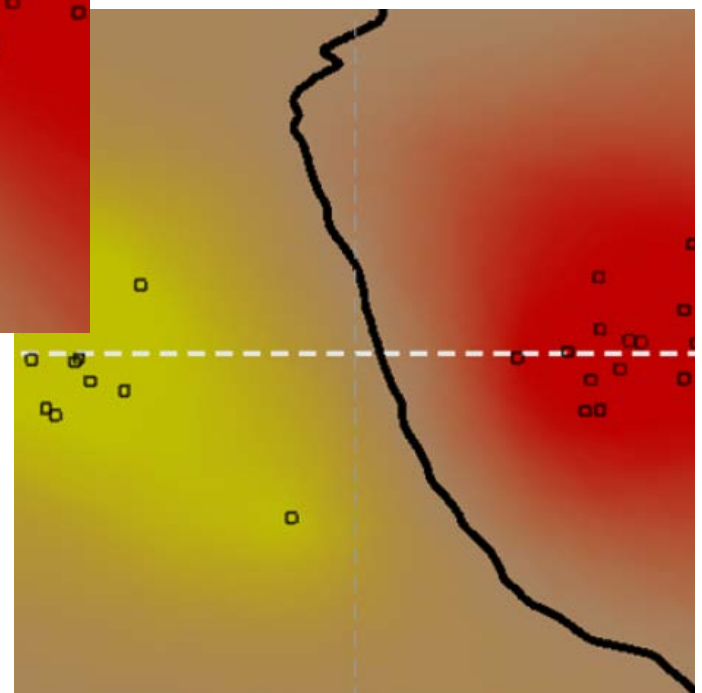
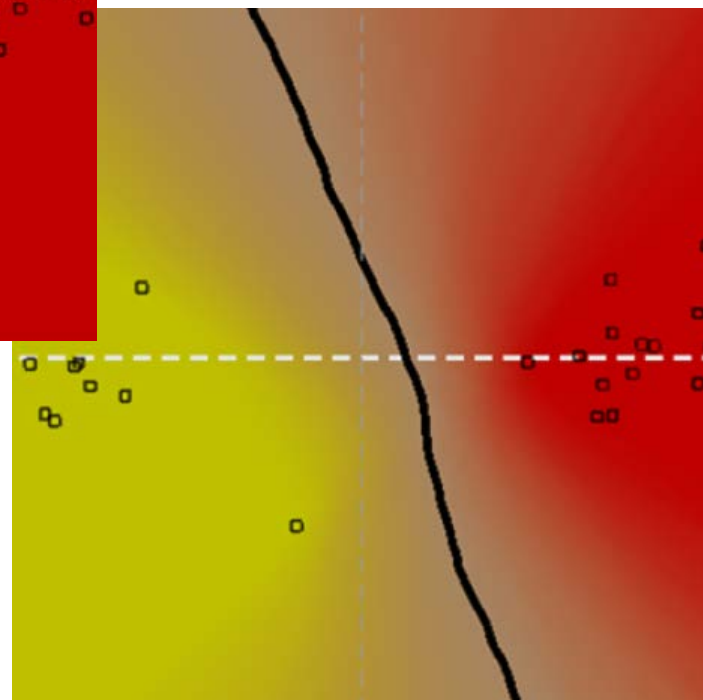
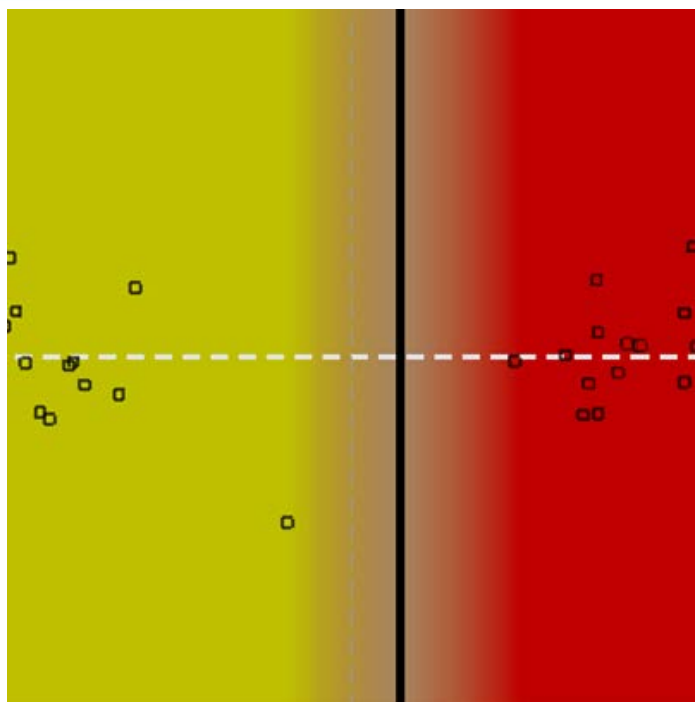
x_2  x_1 x_2  x_1 

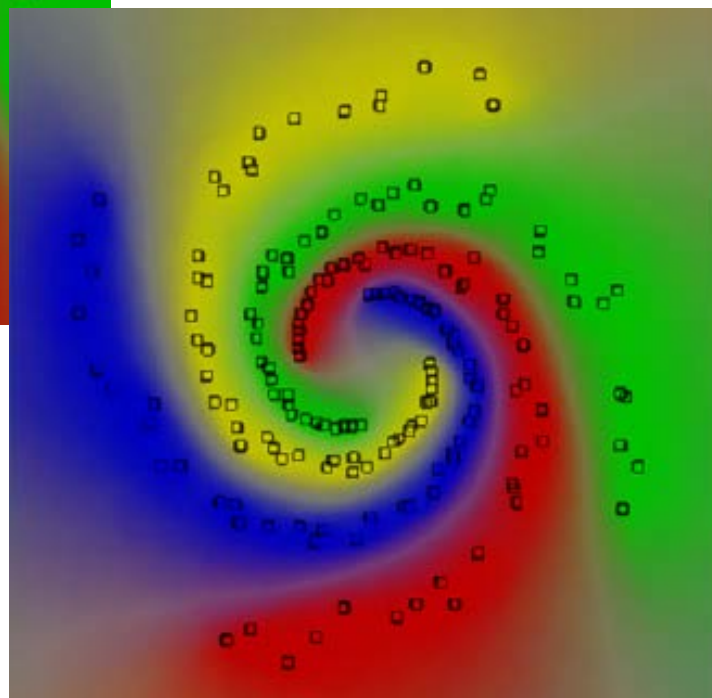
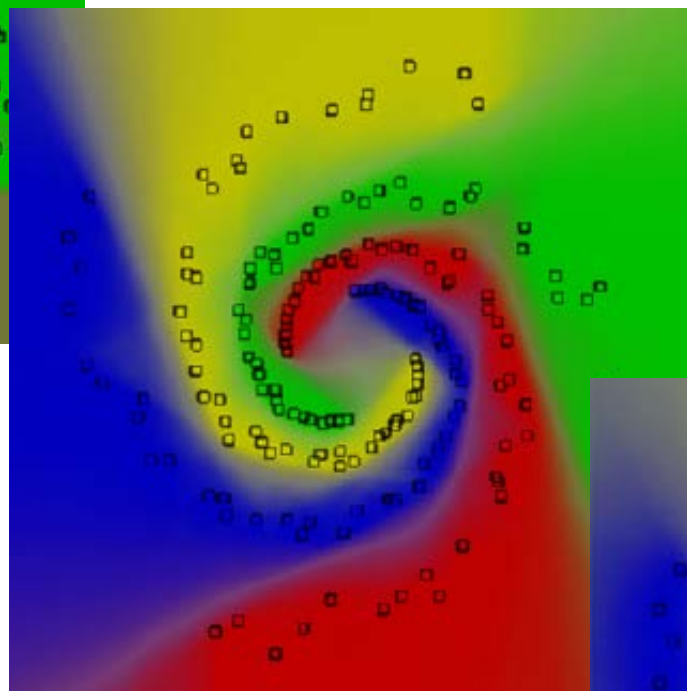
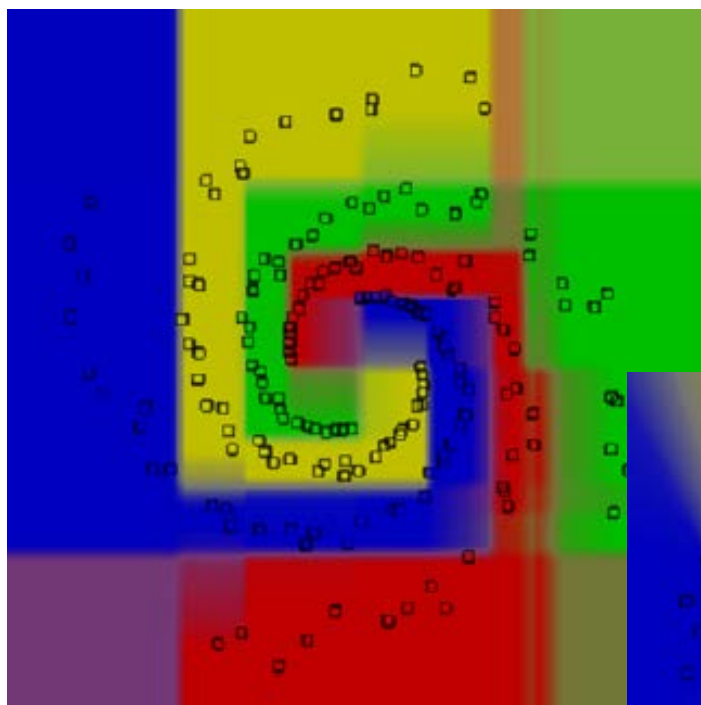


Typy náhodných lesů

1. Diferenciace stromů učením různých podmnožin dat
2. Diferenciace stromů zahrnutím různých množin parametrů



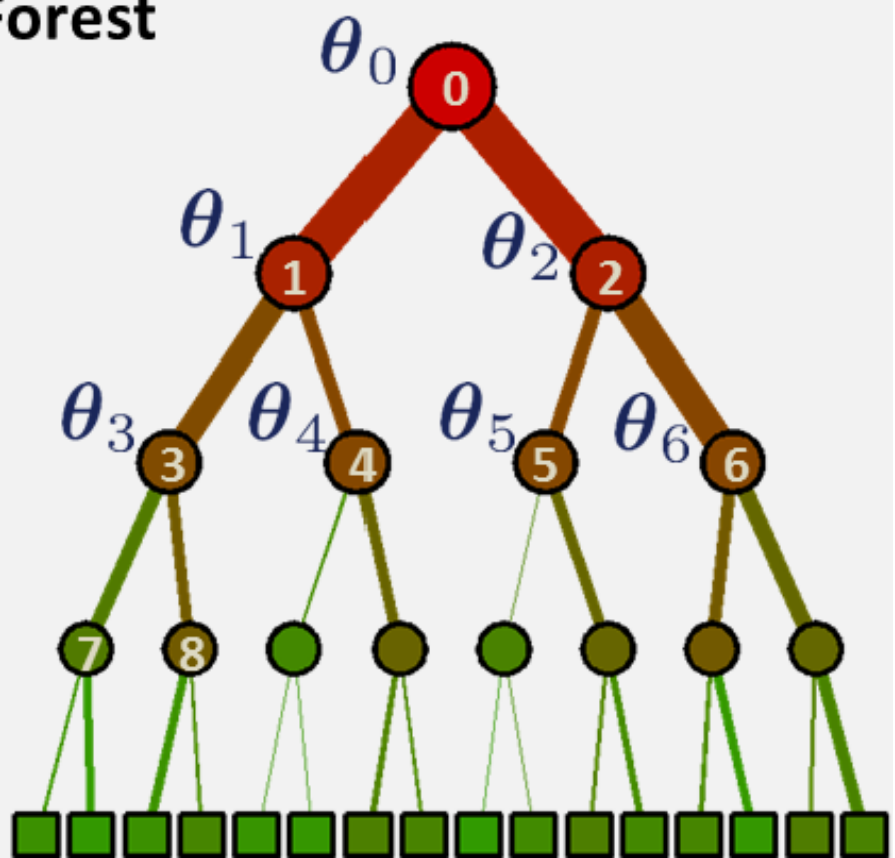




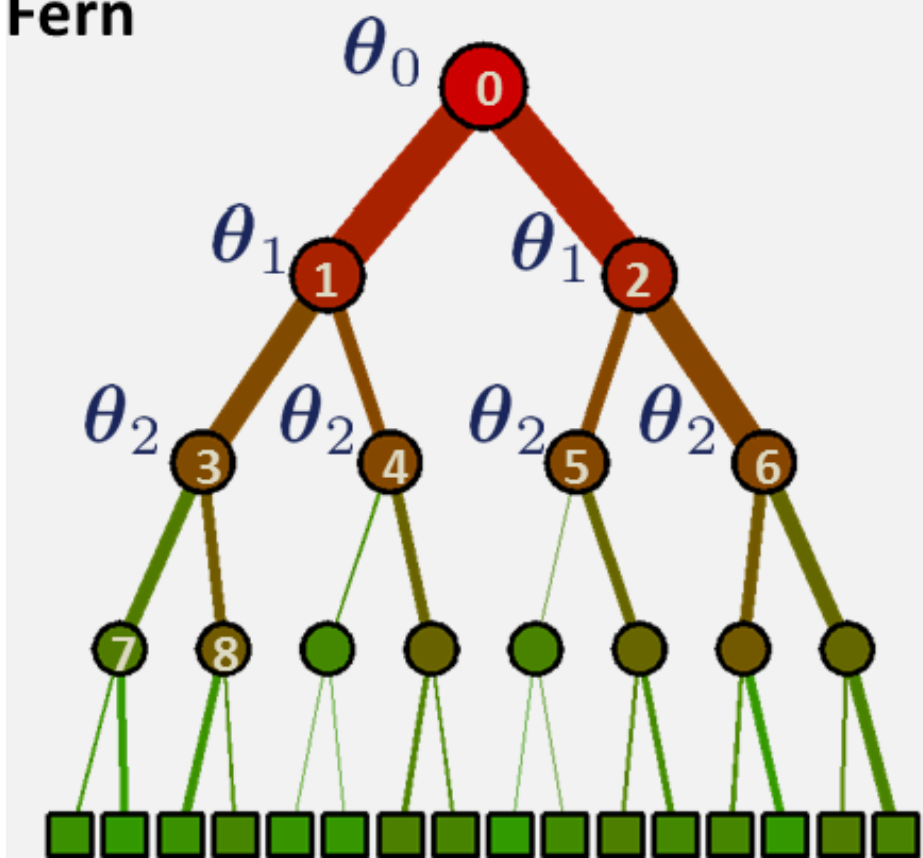
Typy náhodných lesů

1. Diferenciace stromů zahrnutím různých množin parametrů
 2. Diferenciace stromů učením různých podmnožin dat
 3. Diferenciace stromů kombinací 1. a 2.
- ◆ Kapradí – les specifických stromů: uzlům stejné úrovně odpovídají stejné proměnné + hodnoty parametrů

Forest

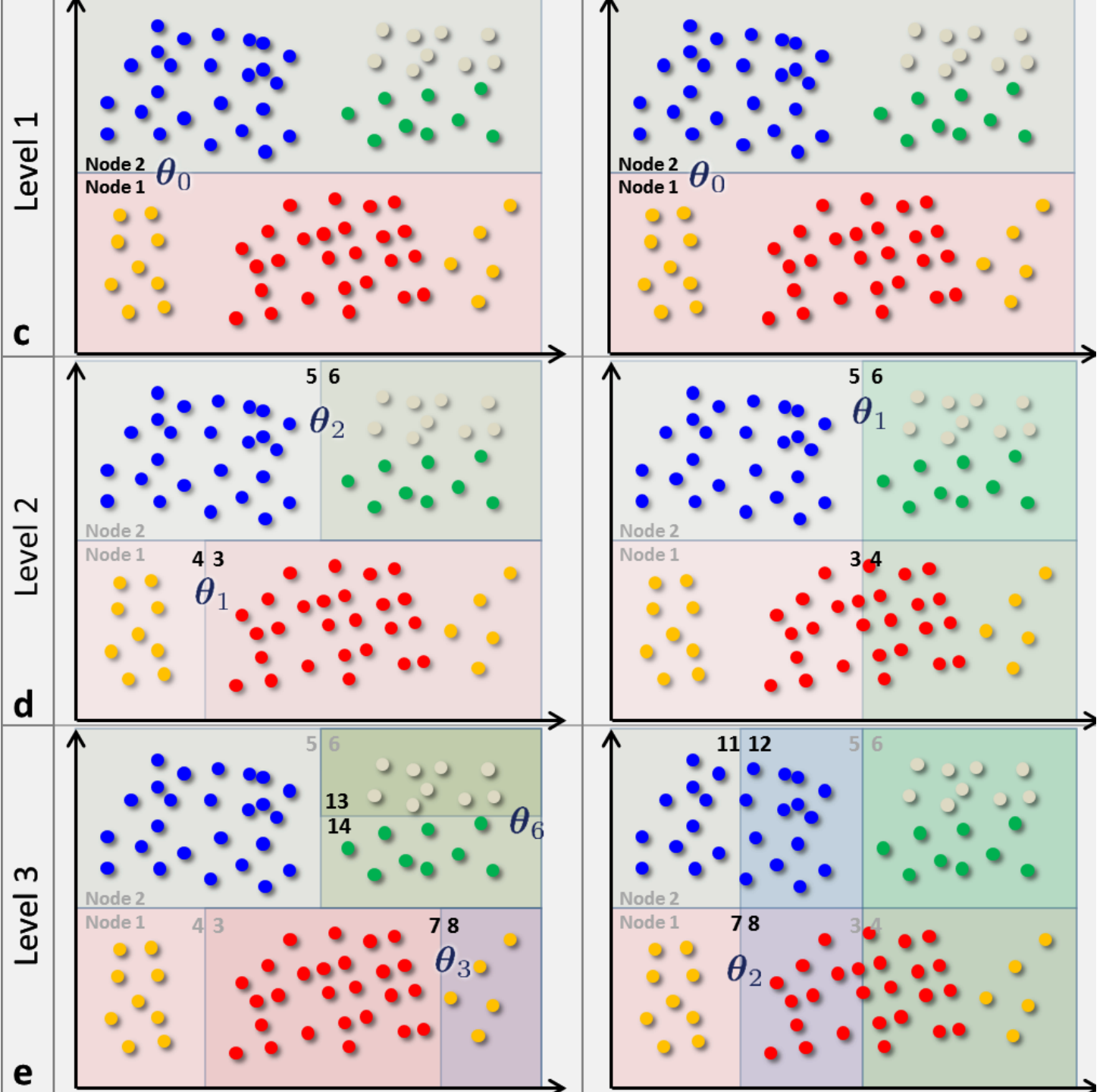


Fern



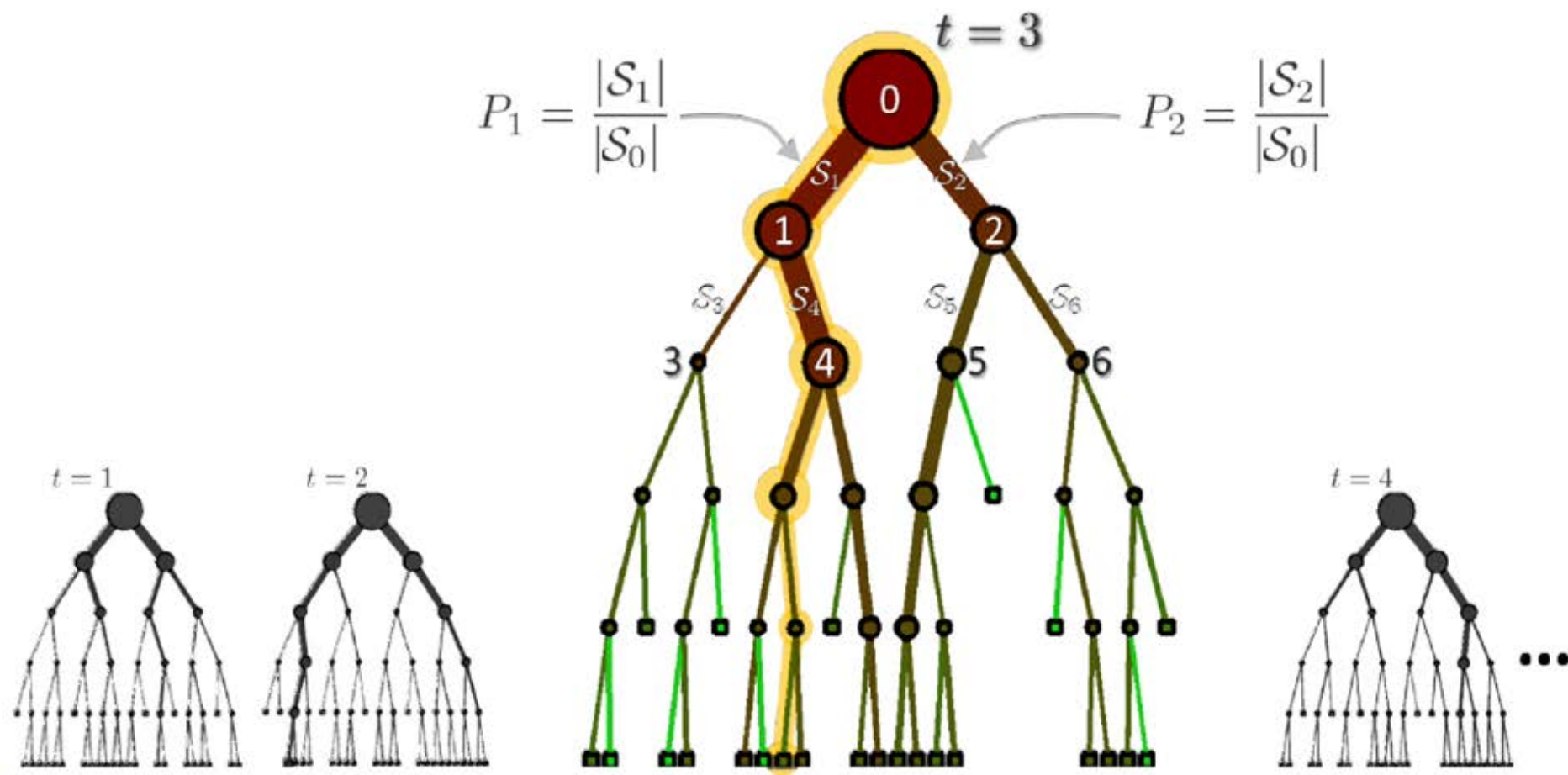
Typy náhodných lesů

1. Diferenciace stromů zahrnutím různých množin parametrů
 2. Diferenciace stromů učením různých podmnožin dat
 3. Diferenciace stromů kombinací 1. a 2.
- ◆ Kapradí – les specifických stromů: uzlům stejné úrovně odpovídají stejné proměnné + hodnoty parametrů
- → rovnoměrnější rozdělení vstupního prostoru (někdy +, jindy -)



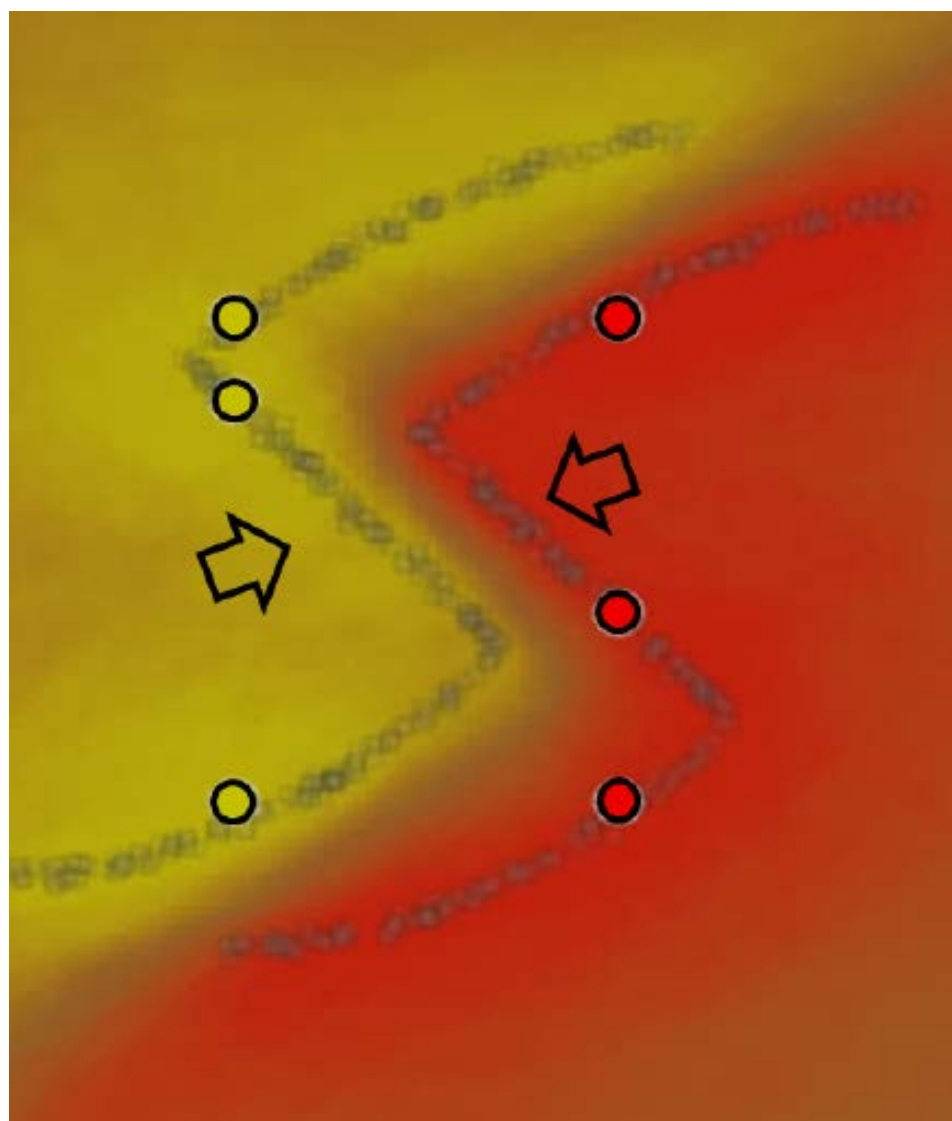
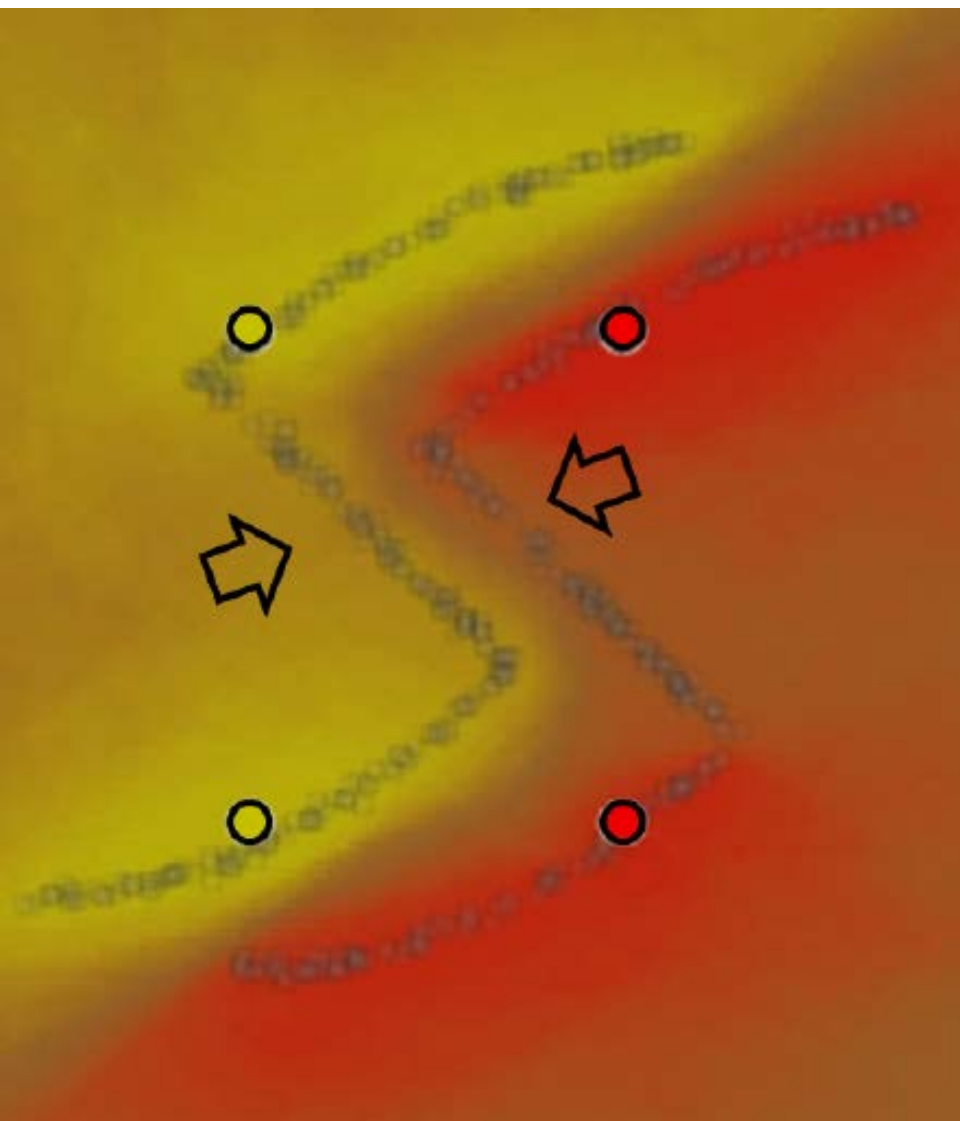
Učení náhodných lesů

- ◆ Každý strom zvlášť (informační zisk | Gini)
 - pouze parametry specificky pro něj zahrnuté
 - pouze na jemu vybrané podmnožině dat
- ◆ Výsledné pravděpodobnosti klasifikace lesem: \sum všechny stromy



Učení náhodných lesů

- ◆ Každý strom zvlášť (informační zisk | Gini)
 - pouze parametry specificky pro něj zahrnuté
 - pouze na jemu vybrané podmnožině dat
- ◆ Výsledné pravděpodobnosti klasifikace lesem: \sum všechny stromy
- ◆ Aktivní učení: drahé zjištění správné klasifikace
 - po SVM 2. nejčastější případ, kdy používáno



Náhodné lesy při doporučování

- ◆ Klasifikování uživatelů do různých demografických profilů
 - *demografický profil*: kombinace věku, pohlaví, vzdělání,...
 - motivace: profilům odpovídají *specifické strategie doporučování*
- ◆ *Učení*: dotázaním některých uživatelů (např. registrovaných)
 - neriskujeme odrazení ostatních vyptáváním – zdržující, vlezlé
- ◆ Hlavní důvod lesů: diskrétní + spojité vstupy