

Internet a klasifikační metody, 1. přednáška

# Tři důležité internetové aplikace klasifikačních metod

volitelný předmět pro magisterské studium

Martin Holeňa



# O čem to bude?

## ◆ Filtrace spamu

- klasifikace spamu na základě obsahu + metainformací

## ◆ Doporučovací systémy

- klasifikace při obsahovém filtrování + kolaborativním filtrování

## ◆ Systémy pro odhalování hrozeb v síti

- klasifikace přítomnosti + druhu anomálního chování sítě

# Filtrace spamu

- ◆ Oddělování spamu od žádoucí pošty (hamu)
- ◆ Filtrace se provádí na několika místech:
  - *filtrační servery* předřazené : spam typicky nedoručen
  - servery *příchozí* pošty } spam doručen, ale
  - poštovní *klienti* } explicitně označen | oddělen

Gmail ▾

SCHREIBEN

Posteingang (2.377)

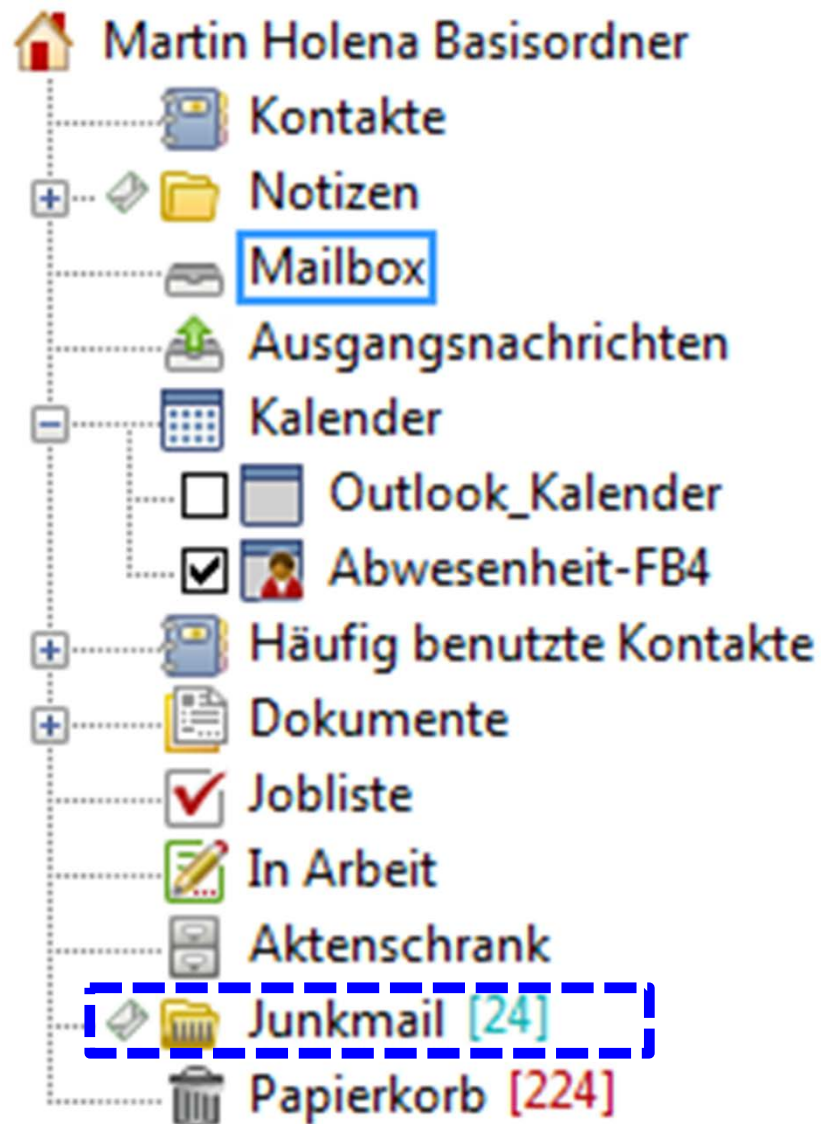
Wichtig

Gesendet

Entwürfe (1)

Spam (323)

Papierkorb



# Filtrace spamu

- ◆ Oddělování spamu od žádoucí pošty (hamu)
- ◆ Filtrace se provádí na několika místech:
  - *filtrační servery* předřazené : spam typicky nedoručen
  - servery *příchozí* pošty } spam doručen, ale
  - poštovní *klienti* } explicitně označen | oddělen
  - servery odchozí pošty: *odchozí (egresní)* filtrace

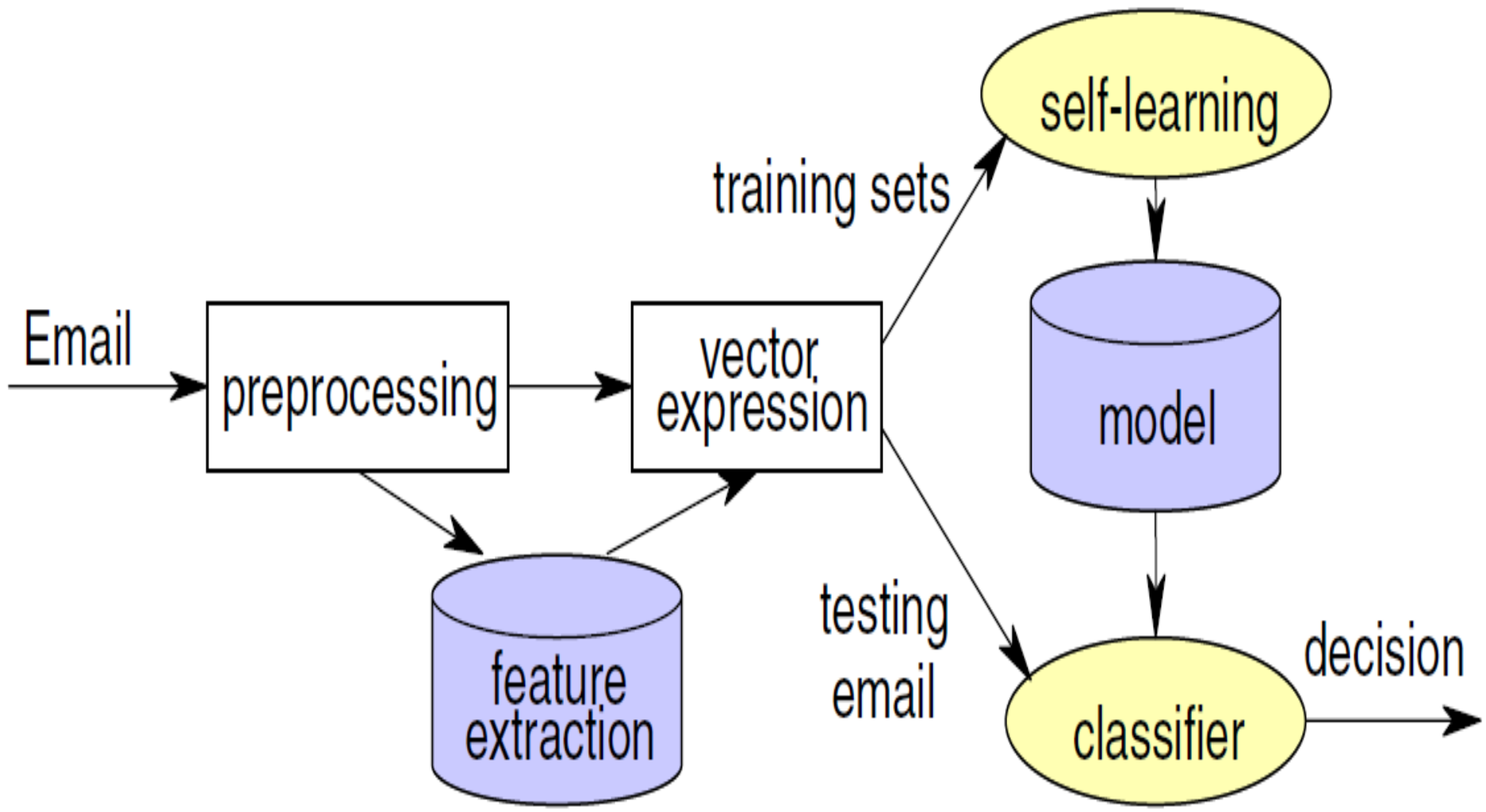
# Proč je filtrace spamu klasifikace?

- ◆ Rozpoznání spamu závisí na mnoha příznacích:
  - obsah zprávy: slova, kombinace slov, obrázky
  - postoj adresáta zprávy k jejímu obsahu
  - atributy hlavičky: předmět, odesílatel, podvržené atributy
- ◆ *{kombinace hodnot jednotlivých příznaků} → {spam, ham}*
  - každé zobrazení {kombinace} → {třídy} je klasifikátor



# Průběh filtrace spamu







# Průběh filtrace spamu

1. Extrakce *hodnot příznaků* ze zprávy
  - použité příznaky určují typ klasifikace:  
podle obsahu zprávy, podle metainformací, hybridní
2. *Učení klasifikátoru* = hledání parametrů jeho modelu
  - každý typ klasifikátoru odpovídá specifickému modelu
3. *Použití* naučeného *klasifikátoru* na nové zprávy

**Unstructured set of tokens: *header***

from,mary,example,  
com,to,mike,org,  
received,...

**Selected fields  
of the header**

IP<sub>1</sub> = [xxx.xxx.xxx.xxx]  
IP<sub>2</sub> = [yyy.yyy.yyy.yyy]  
...

**Unstructured set  
of tokens : *all***

from, mary,example,  
com, to, mike, org,  
received,...  
dear,i,would,like ...

From: <mary@example.com>  
 To: <mike@example.org>  
 Received: from [xxx.xxx.xxx.xxx] by ...  
 Received: from [yyy.yyy.yyy.yyy] by ...  
 ...

---

 Dear Mike!  
 I would like to  
 congratulate  
 you with ...

**General characteristics**

Size = 2, 411  
NumberOfAttachments = 0  
...

**Unstructured set of  
tokens : *body***

dear,mike,i,would,  
like,to,congratulate,  
...

**Graphical elements**



**Body as a text in a  
natural language**

Dear Mike!  
I would like to  
congratulate you  
with ...

# Filtrace podle obsahu

- ◆ Používané příznaky: *slova*, *kombinace* slov, *obrázky*
- ◆ Učení klasifikátorů filtrujících na základě obsahu:
  - Možnost *doučování* zpětnou vazbou od adresátů
    1. *globální* klasifikátor – zpětnou vazbou všech
    2. *lokální* klasifikátory – zpětnou vazbou jednotlivých adresátů(doučování majitelem sexshopu: „levná viagra“ = ham)

# Filtrace podle metainformací

- ◆ Metainformace (= informace o zprávě) z *hlavičky*
  - doména, kódování (azbuka), podvržená adresa odesilatele
- ◆ Metainformace z *průběhu posílání* zpráv
  - *jedné*: nedodržení SMTP ukončení spojení (QUIT),...
  - *více*: stejná/podobná zpráva mnoha adresátům
- ◆ Kombinování s obsahem zprávy → *hybridní* filtry

# Příklady spamových filtrů

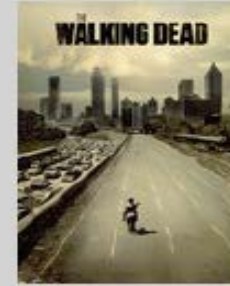
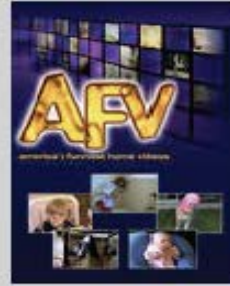
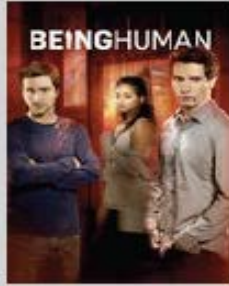
- ◆ Součást běžných *mailových klientů* + *webmailových serverů*
  - Outlook, Groupwise, Gmail, Email.cz (Seznam)
- ◆ *SpamBayes* – pythonovský program předřazený klientovi
  - Klasifikuje do 3 tříd: spam, ham, nejistota
- ◆ *Další* filtry ve formě samostatných programů:
  - DSPAM, SpamAssassin, Bogofilter, ASSP, Babletext

# Doporučovací systémy

- ◆ Doporučují uživateli, co ho pravděpodobně zaujme:  
zboží, literatura, filmy, rekreace, jiní uživatelé

# Thousands of movies and TV episodes including these:

## New Arrivals in TV



## TV Drama



## TV Comedy



## Children & Family



Netflix

# Personen, die Sie vielleicht auch kennen

Sehen Sie Personen aus verschiedenen Bereichen Ihres Berufslebens

Univerzita  
Karlova v  
Praze

České  
vysoké  
učení



FIT CTU in  
Prague

Vysoká  
škola  
ekonomická



**Vera Kurkova** 2

scientist at Institute of Computer Science  
Academy of Sciences of the Czech Republic  
Tschechische Republik

[+ Vernetzen](#)

4 gemeinsame Kontakte



**Ladislav Beneš** ×

chief IT at UI AV  
Tschechische Republik

[+ Vernetzen](#)



**Martin Štekl** 3

PHP Developer at Designo Creative s.r.o.  
Tschechische Republik

[+ Vernetzen](#)



**Petr Hajek** 2

vedecký pracovník ve společnosti Ustav  
informatiky AV  
Prag, Tschechische Republik

[+ Vernetzen](#)

1 gemeinsamer Kontakt



**Jiri Palek** 3

PhD student at CTU, FNSPE; OSVČ  
Bezirk Rakonitz, Tschechische Republik

[+ Vernetzen](#)



**Tomáš Kalvoda**

teacher, researcher at FIT ČVUT  
Prag, Tschechische Republik

[+ Vernetzen](#)



**Lukáš Hošek** 3

Programmer, Bohema Interactive Simulations /  
PhD candidate, MFF UK  
Prag, Tschechische Republik

[+ Vernetzen](#)



**Jitka Hodalová** 2

secondary school teacher (SOŠ Emila Holuba,  
s.r.o.)  
Bezirk Brünn-Stadt, Tschechische Republik

[+ Vernetzen](#)

1 gemeinsamer Kontakt



# Doporučovací systémy

- ◆ Doporučují uživateli, co ho pravděpodobně zaujme:  
zboží, literatura, filmy, rekreace, jiní uživatelé
- ◆ Výstup: *1. buď* jen doporučený/é *objekt/y* zájmu  
*2. nebo* seznam s *odhadem* pravděpodobné *zajímavosti*
- ◆ Často *vysvětlují*, proč objekt zájmu doporučují  
→ vyšší srozumitelnost doporučení pro uživatele

## Frequently Bought Together



Price For All Three: **\$181.05**



Add all three to Cart

Add all three to Wish List

Show availability and shipping details

- This item:** Gaussian Processes for Machine Learning (Adaptive Computation and Mac Edward Rasmussen Hardcover **\$35.97**
- Pattern Recognition and Machine Learning (Information Science and Statistics) by Ch **\$61.26**
- Machine Learning: A Probabilistic Perspective (Adaptive Computation and Machine Le Murphy Hardcover **\$83.82**

## Customers Who Bought This Item Also Bought



Machine Learning: A Probabilistic ...

> Kevin P. Murphy

★★★★★ (5)

Hardcover

**\$83.82**



Pattern Recognition and Machine Learning ...

> Christopher M. Bishop

★★★★★ (67)

Hardcover

**\$61.26**



Probabilistic Graphical Models: Principles and ...

> Daphne Koller

★★★★★ (16)

Hardcover

**\$89.89**


# Proč je doporučování klasifikace?

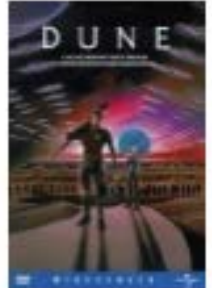
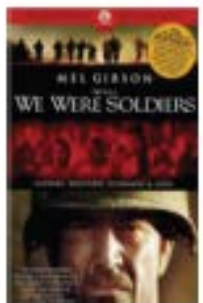
- ◆ Rozhodnutí zda + co doporučit závisí na:
  - *příznamech* popisujících *obsah* objektu zájmu
  - příznamech *chování uživatele*, kterému je doporučováno
  - příznamech *chování ostatních* uživatelů
- ◆ → 2 typy rozhodování zda + co doporučit, podle příznaků: 1. obsahové filtrování, 2. kolaborativní filtrování

# Obsahové filtrování

- ◆ Doporučovaný objekt vybrán podle *podobnosti* s
  - objekty vyhodnoceými jako *zajímaví* uživatele *dosud*:  
*explicitně* (hodnocení) | *implicitně* (koupě, stažení, streamování)
  - někdy: explicitním stanovením určitých příznaků uživatelem
- ◆ *Příklady*: Internet Movie Database — doporučování videí,  
Pandora Radio — skladby „podobné“ původně zvoleným

# Kolaborativní filtrování

- ◆ Princip: uživatelé s *podobnou historií chování* se často *zajímají* o *podobné* objekty
  - historie chování: zadávané dotazy, prolézané stránky, projevený zájem o stejné objekty: explicitně, implicitně
- ◆ Problém náběhu systému: o  málo dat
- ◆ *Hybridní* systémy: obsahové i kolaborativní filtrování



# Kolaborativní a hybridní příklady

- ◆ *Amazon*: 2 způsoby využití informace kupující → zboží
- ◆ *Last.fm*: skladby dle podobné historie poslouchání
- ◆ *LinkedIn, Facebook*: doporučují kontakty + členské skupiny

# Personen, die Sie vielleicht auch kennen

Sehen Sie Personen aus verschiedenen Bereichen Ihres Berufslebens

Univerzita  
Karlova v  
Praze

České  
vysoké  
učení



FIT CTU in  
Prague

Vysoká  
škola  
ekonomická



**Vera Kurkova** 2

scientist at Institute of Computer Science  
Academy of Sciences of the Czech Republic  
Tschechische Republik

[+ Vernetzen](#)

4 gemeinsame Kontakte



**Ladislav Beneš** ×

chief IT at UI AV  
Tschechische Republik

[+ Vernetzen](#)



**Martin Štekl** 3

PHP Developer at Designo Creative s.r.o.  
Tschechische Republik

[+ Vernetzen](#)



**Petr Hajek** 2

vedecký pracovník ve společnosti Ustav  
informatiky AV  
Prag, Tschechische Republik

[+ Vernetzen](#)

1 gemeinsamer Kontakt



**Jiri Palek** 3

PhD student at CTU, FNSPE; OSVČ  
Bezirk Rakonitz, Tschechische Republik

[+ Vernetzen](#)



**Tomáš Kalvoda**

teacher, researcher at FIT ČVUT  
Prag, Tschechische Republik

[+ Vernetzen](#)



**Lukáš Hošek** 3

Programmer, Bohema Interactive Simulations /  
PhD candidate, MFF UK  
Prag, Tschechische Republik

[+ Vernetzen](#)



**Jitka Hodalová** 2

secondary school teacher (SOŠ Emila Holuba,  
s.r.o.)  
Bezirk Brunn-Stadt, Tschechische Republik

[+ Vernetzen](#)

1 gemeinsamer Kontakt

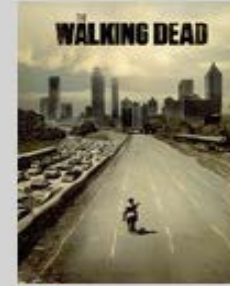
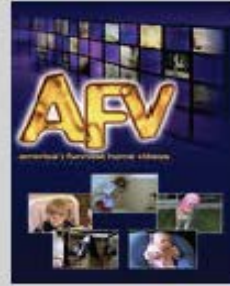
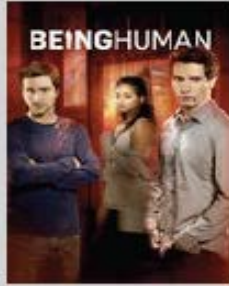


# Kolaborativní a hybridní příklady

- ◆ *Amazon*: 2 způsoby využití informace kupující → zboží
- ◆ *Last.fm*: skladby dle podobné historie poslouchání
- ◆ *LinkedIn, Facebook*: doporučují kontakty + členské skupiny
- ◆ *hybridní systémy* Netflix, See This Next:
  - videa podobná těm, co uživatel chválil
  - + dle uživatelské podobnosti historie půjčování + hledání

# Thousands of movies and TV episodes including these:

## New Arrivals in TV



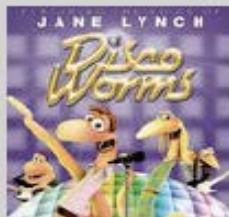
## TV Drama



## TV Comedy



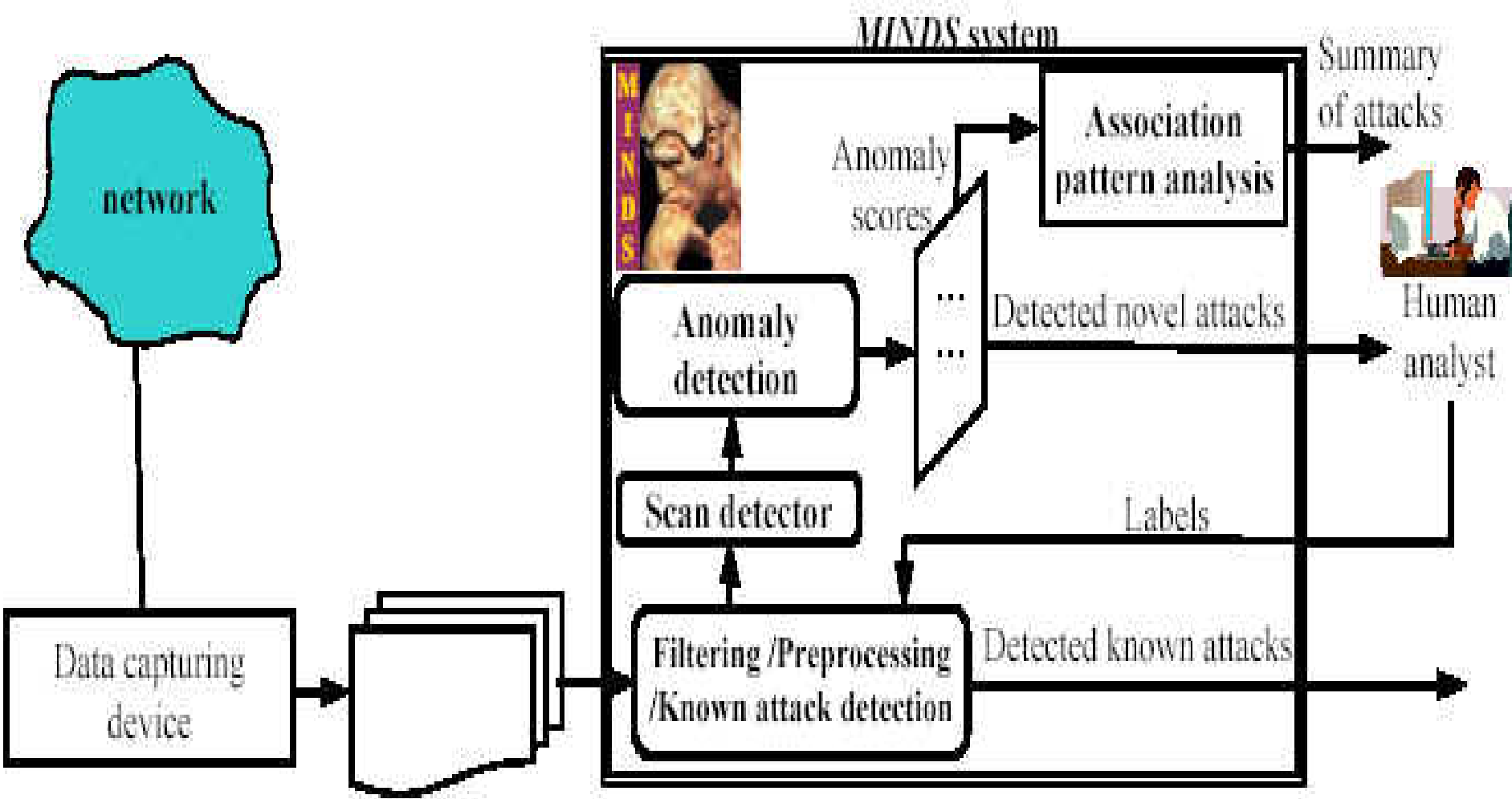
## Children & Family



Netflix

# Systemy pro odhalení hrozeb v síti

- ◆ Systemy odhalující v síti pokusy o:
  - neoprávněný síťový *přístup* (remote to local)
  - neoprávněná *administrátorská práva* (local to root)
  - *zabránění využívání zdrojů* (denial of service)
  - shromažďování informací o slabínách zabezpečení sítě
- ◆ Příklad: Minnesota Intrusion Detection System (MINDS)



# Proč je odhalování hrozeb klasifikace?

- ◆ K odhalení přítomnosti hrozby v síti

je třeba znát hodnoty řady příznaků:

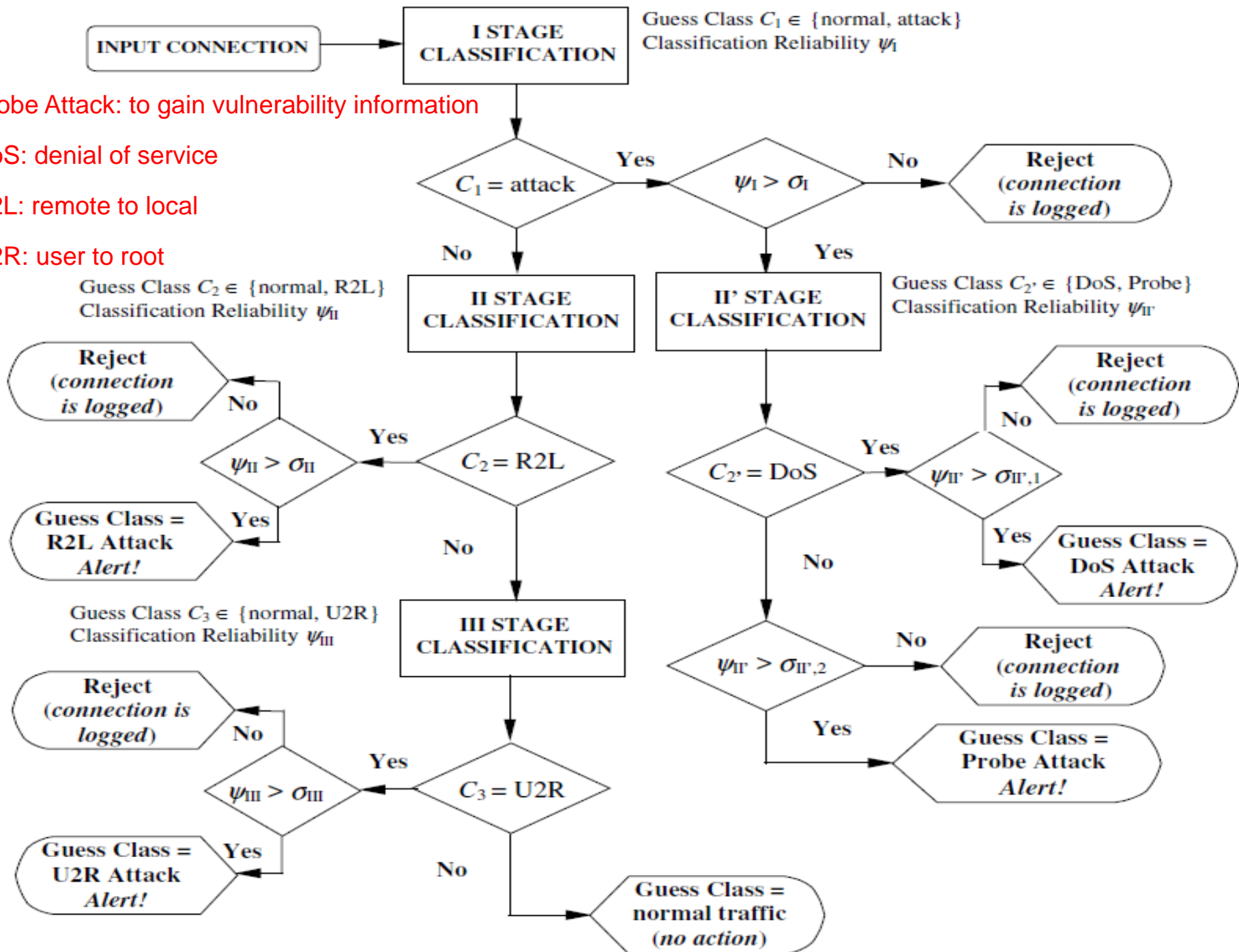
- zdrojové + cílové adresy, protokoly, počty paketů, ...

- ◆ {kombinace hodnot příznaků} → {normál, anomálie}

{normál, R2L, L2R, DoS, ...}

- ◆ *Signatury: známé*, s pravděpodobností  $\cong 100\%$

Probe Attack: to gain vulnerability information  
 DoS: denial of service  
 R2L: remote to local  
 U2R: user to root



# Souvislost s data-miningem

- ◆ Příznaky používané systémem pro odhalení hrozeb ke klasifikaci vytváří tzv. auditová data
- ◆ Objevování nových znalostí v  $\nearrow$  = data mining
  - česky: dobývání znalostí, vytěžování dat, datokopectví
  - popisné statistiky, modely časových řad, shlukování, ...
- ◆ Systém ADAM (Audit Data Analysis & Mining)

# Co vám tento předmět chce dát?

- ◆ Povědomí o velkém spektru klasifikačních metod
  - o jejich *zajímavosti, důmyslnosti, užitku* pro internet
- ◆ Základní vybavení pro jejich tvůrčí používání
  - s *porozuměním*, vhladem pod povrch, samostatně
  - nejenom na problémy řešené jimi *dnes*
  - i problémy, které potkáte za *10–20 roků*



# O čem si proto budeme povídat?

1. Základní *koncepty* týkající se klasifikace
  2. Hlavní typy klasifikačních *metod*
  3. *Přesnost* klasifikace nově příchozích dat
  4. *Srozumitelnost* výsledku klasifikace pro uživatele
  5. Spojování klasifikátorů do *týmů*
- ◆ Dnešní *3 aplikace* všechno ilustrují, nikoliv vymezují

# Cvičení: peer přednášky

- ◆ 6 cvičení střídavě s přednáškami
  - ve výjovém prostředí Matlab (seznámení: 1. cvičení)
- ◆ 2 stupně obtížnosti procvičované látky:
  - jednoduché příklady procvičující probírané téma
  - 1 náročnější komplexní úkol jako semestrální práce
  - při cvičeních i konzultace vybraných semestrálek