



# ANALÝZA A KLASIFIKACE DAT



prof. Ing. Jiří Holčík, CSc.



INVESTICE DO ROZVOJE VZDĚLÁVÁNÍ

# LITERATURA

---

- ☑ Holčík, J.: přednáškové prezentace
- ☑ Holčík, J.: Analýza a klasifikace signálů. [Učební texty VŠ], Brno, FE VUT 1992.

# LITERATURA

- ✓ Duda,R.O., Hart,P., Stork,D.G. Pattern Classification. New York, John Wiley & Sons 2001
- ✓ Theodoridis S., Koutroumbas K., Pattern Recognition. Amsterdam, Elsevier 2009
- ✓ McLachlan,G.J.: Discriminant Analysis and Statistical Pattern Recognition. J.Wiley&Sons, Hoboken 2004
- ✓ Webb,A.: Statistical Pattern Recognition. J.Wiley&Sons, Chichester 2002
- ✓ Meloun, M., Militký,J.: Statistická analýza experimentálních dat. Praha, Academia 2004.



# 0. ČEM TO BUDE?



# ANOTACE

Předmět poskytne informaci o základních metodách a algoritmech pro výběr popisu, hodnocení a klasifikaci biomedicínských dat. Zabývá se základním tříděním klasifikačních přístupů – příznakové a strukturální a uvádí principy obou přístupů. Dále se zabývá podrobně zejména metodami příznakovými. Klasifikace podle diskriminačních funkcí (princip a stanovení diskriminačních funkcí na základě statistických vlastností množiny obrazů) a minimální vzdálenosti. Sekvenční klasifikace. Volba a výběr příznaků. Selektce a extrakce příznaků – analýza hlavních a nezávislých komponent, faktorová analýza. Učení klasifikátorů. Shlukování – podobnost mezi obrazy, podobnost mezi shluky, metody shlukování. Klasifikace pomocí neuronových sítí.

# ANOTACE

Předmět poskytne informaci o základních metodách a algoritmech pro výběr popisu, hodnocení a klasifikaci biomedicínských dat. Zabývá se základním tříděním klasifikačních přístupů – příznakové a strukturální a uvádí principy obou přístupů. Dále se zabývá podrobně zejména metodami příznakovými. Klasifikace podle diskriminačních funkcí (princip a stanovení diskriminačních funkcí na základě statistických vlastností množiny obrazů) a minimální vzdálenosti. Sekvenční klasifikace. Volba a výběr příznaků. Selektce a extrakce příznaků – analýza hlavních a nezávislých komponent, faktorová analýza. Učení klasifikátorů. Shlukování – podobnost mezi obrazy, podobnost mezi shluky, metody shlukování. Klasifikace pomocí neuronových sítí.

# OSNOVA

- ☑ Klasifikace dat – základní terminologie. Klasifikace vs. diskriminační analýza vs. predikce. Klasifikace vs. regrese. Třídění klasifikačních algoritmů - klasifikace pomocí minimální vzdálenosti, pomocí ztotožnění s etalony, pomocí diskriminačních funkcí (lineární, nelineární), pomocí definice hranic mezi jednotlivými třídami.
- ☑ Parametrické vs. neparametrické přístupy. Učení s učitelem, bez učitele, s nedokonalým učitelem.
- ☑ Strukturální popis a klasifikace. Primitiva a relace, hierarchický a nehierarchický popis, reprezentace klasifikačních tříd pomocí gramatiky, automatu. Strukturální metriky.
- ☑ Příznakové metody. – Příznak, znak, diskriminátor, prediktor. Klasifikace podle minimální vzdálenosti – metrika, funkce podobnosti, vzdálenost mezi obrazy, vzdálenost mezi obrazem a množinou obrazů. Příklady metrik – deterministické, pravděpodobnostní. Příklady funkcí podobnosti - asociační koeficienty, korelační koeficienty.
- ☑ Příznaková klasifikace podle diskriminačních funkcí – Fisherův algoritmus, Bayesův klasifikátor. Stanovení diskriminačních funkcí na základě statistických vlastností množiny obrazů.

# OSNOVA

- ☑ Příznaková klasifikace podle diskriminačních funkcí – Fisherův algoritmus, Bayesův klasifikátor. Stanovení diskriminačních funkcí na základě statistických vlastností množiny obrazů.
- ☑ Lineární diskriminační funkce – dichotomický a multikategoriální problém, zobecněné lineární diskriminační funkce. Lineárně separabilní a neseperabilní případy. Logistická diskriminace.
- ☑ Kontextová klasifikace – Bayesův klasifikátor, Markovovy modely, Viterbiho klasifikátor, skryté Markovovy modely,
- ☑ Volba a výběr příznaků. Selektce a extrakce (generování) příznaků, Transformace dat a redukce dimenzionality. Ordinační metody. Kritéria a algoritmy selektce příznaků.
- ☑ Faktorová analýza – princip, důsledky.
- ☑ Analýza komponent. Analýza hlavních komponent – princip, důsledky.
- ☑ Analýza nelineárních komponent – princip, důsledky. Analýza nezávislých komponent – princip, důsledky.
- ☑ Sekvenční klasifikace. Princip, Waldovo a Reedovo kritérium, jejich modifikované varianty.



# UKONČENÍ PŘEDMĚTU

Požadavky:

☑ ústní zkouška

→ dvě části:

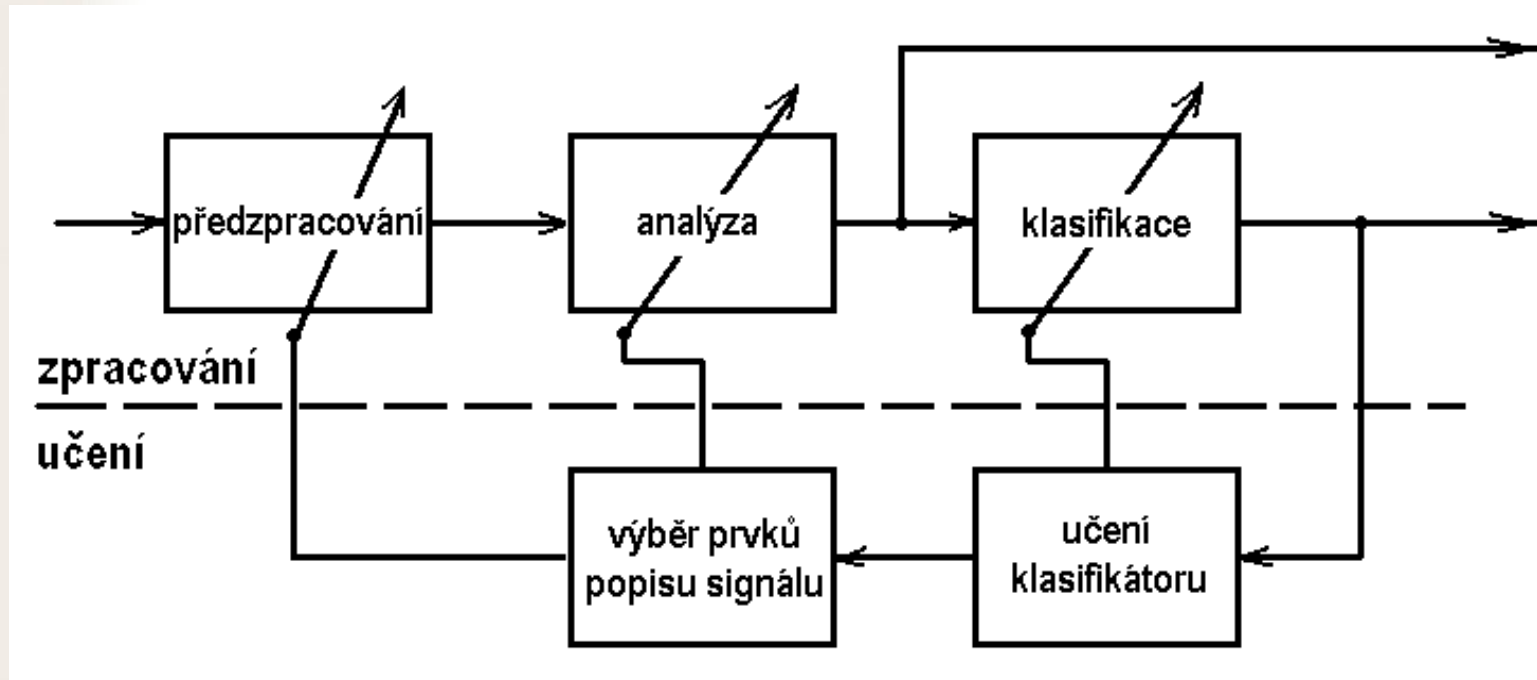
- ☐ učená rozprava o některém z témat, která budou náplní předmětu;
- ☐ diskuze nad vyřešeným problémem týkajícím se problematiky klasifikace dat **a používajícím některé z technik, které budou náplní předmětu;**



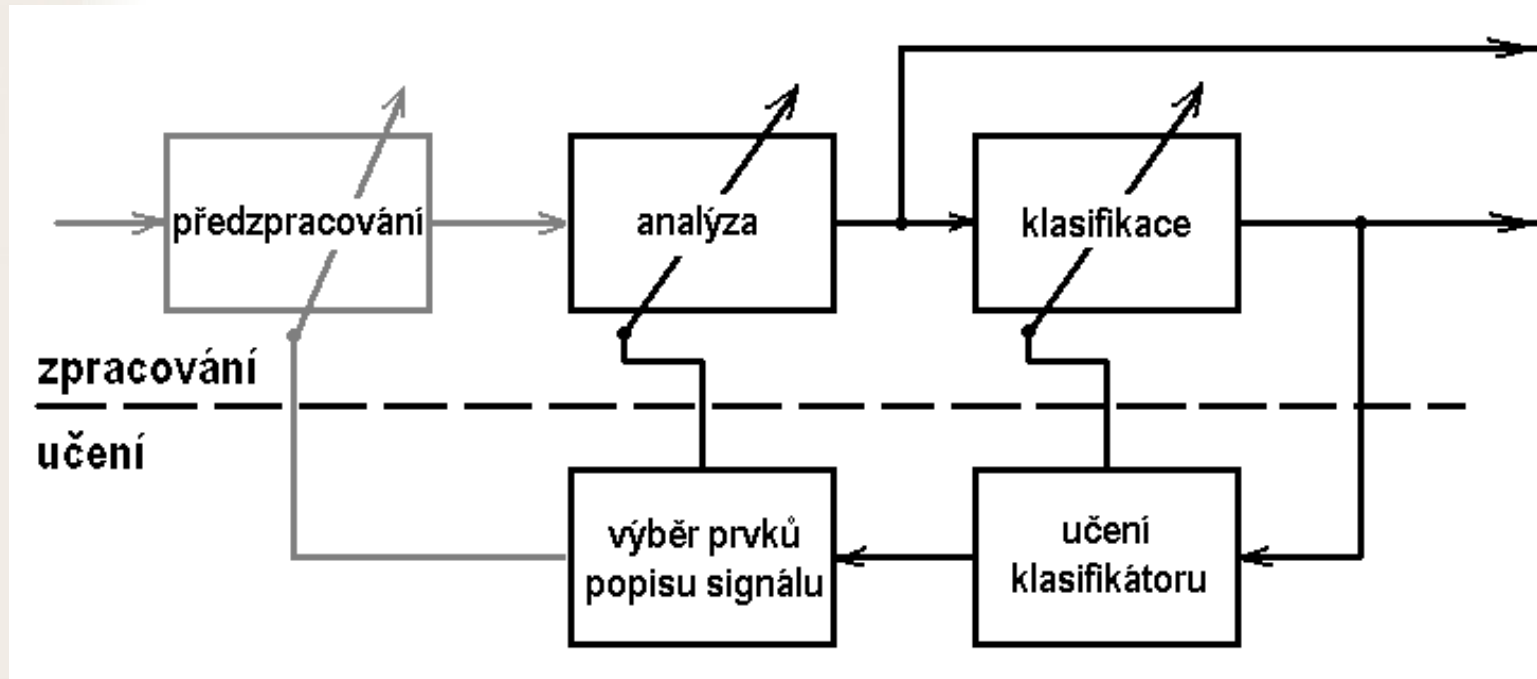
# I. ZAČÍNÁME



# OBEČNÉ SCHÉMA ZPRACOVÁNÍ DAT



# OBECNÉ SCHÉMA ZPRACOVÁNÍ DAT



# OBECNÉ SCHÉMA ZPRACOVÁNÍ DAT

## ZPRACOVÁNÍ

### ☑ předzpracování

- filtrace rušivých složek x zvýraznění užitečných složek signálu;
- rekonstrukce a doplnění chybějících údajů;
- konverze typu dat;
- redukce dat;
- (A/Č převod);

### ☑ analýza dat

- určení hodnot příznaků (reprezentativních parametrů) – pro příznakové klasifikátory;
- nalezení primitiv (charakteristických tvarových segmentů) – strukturální klasifikátory

### ☑ klasifikátor –

- zatřídění do diagnostických kategorií

# OBECNÉ SCHÉMA ZPRACOVÁNÍ DAT

## ZPRACOVÁNÍ

### ☑ **předzpracování**

- filtrace rušivých složek x zvýraznění užitečných složek signálu;
- rekonstrukce a doplnění chybějících údajů;
- konverze typu dat;
- redukce dat;
- (A/Č převod);

### ☑ **analýza dat**

- určení hodnot příznaků (reprezentativních parametrů) – pro příznakové klasifikátory;
- nalezení primitiv (charakteristických tvarových segmentů) – strukturální klasifikátory

### ☑ **klasifikátor –**

- zatřídění do diagnostických kategorií

# OBECNÉ SCHÉMA ZPRACOVÁNÍ DAT

- ☑ **Analýza** (z řečtiny – *rozbor, rozčlenění*) je vědecká metoda založená na dekompozici celku na elementární části. Cílem analýzy je identifikovat podstatné a nutné vlastnosti elementárních částí celku, poznat jejich podstatu a zákonitosti.



# OBECNÉ SCHÉMA ZPRACOVÁNÍ DAT

- ☑ **Analýza** (z řečtiny – *rozbor, rozčlenění*) je vědecká metoda založená na dekompozici celku na elementární části. Cílem analýzy je identifikovat podstatné a nutné vlastnosti elementárních částí celku, poznat jejich podstatu a zákonitosti.
- ☑ **Syntéza** je obecné označení pro proces spojení dvou nebo více částí do jednoho celku. S tímto pojmem se lze setkat v různých spojeních: syntéza obrazu, syntéza řeči, syntéza zvuku, chemická syntéza, jaderná syntéza, termonukleární syntéza, syntéza látek, fotosyntéza, proteosyntéza, biosyntéza, evoluční syntéza.



# ANALÝZA

V bloku analýzy se vytváří formální (abstraktní) popis zpracovávaných dat, který nese **podstatnou** informaci z hlediska kvality rozhodování při klasifikaci. Abstraktní popis se často nazývá **obrazem (pattern)** ⇒ **rozpoznávání obrazů (pattern recognition)**. V datech je vybrána určitá množina elementárních vlastností, příp. jejich elementárních částí a jejich vazeb, jejichž způsob popisu je apriori znám.

# KLASIFIKACE

- ☑ rozumí se rozdělení (konkrétní či teoretické) dané skupiny (množiny) předmětů či jevů na **konečný** počet dílčích skupin (podmnožin), v nichž všechny předměty či jevy mají dostatečně podobné společné vlastnosti. Vlastnosti podle nichž lze klasifikaci zadat či provádět, určují **klasifikační kritéria**. Předměty (jevy), které mají podobnou uvažovanou vlastnost tvoří třídu.

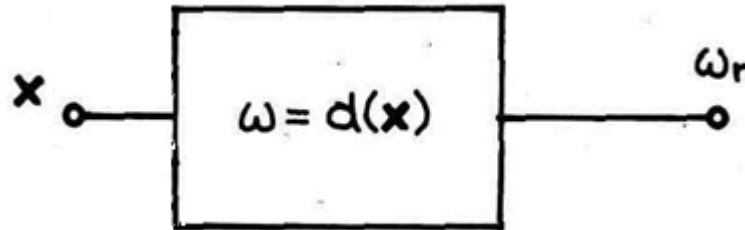
# KLASIFIKÁTOR

- ☑ **Klasifikátor** je stroj (algoritmus,...) s jedním diskretním výstupem, který udává třídu, do které klasifikátor zařadil vstupní reprezentaci dat

$$\omega_r = d(\mathbf{x})$$

$d(\mathbf{x})$  je funkce argumentu  $\mathbf{x}$  představujícího reprezentaci vstupních dat, kterou nazýváme **rozhodovací pravidlo klasifikátoru**;

$\omega_r$  je **identifikátor klasifikační třídy**;  $\omega_r \mid r=1,\dots,R \in \Omega$



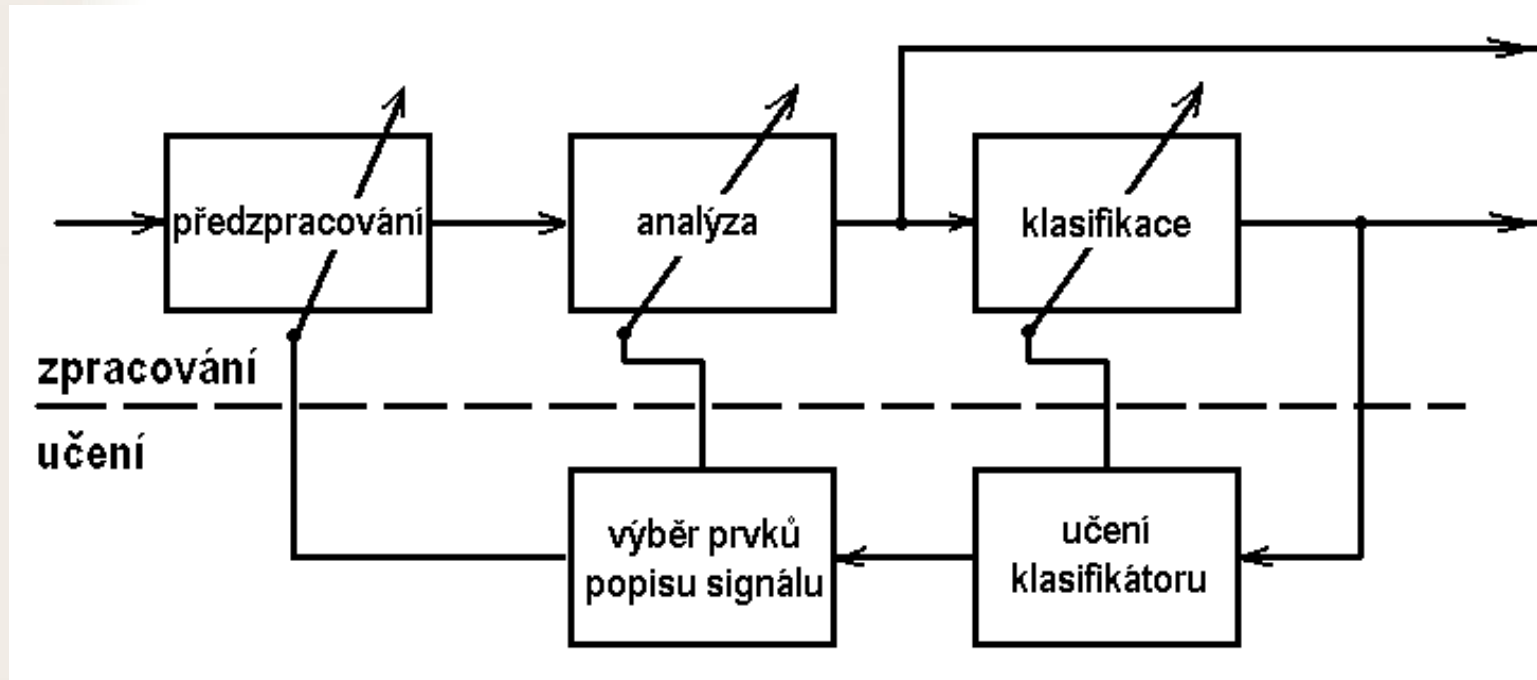
# PRINCIPY KLASIFIKACE

---

# PRINCIPY KLASIFIKACE

- ✓ pomocí **diskriminačních funkcí** – funkcí, které určují míru příslušnosti k dané klasifikační třídě;
- ✓ pomocí **definice hranic** mezi jednotlivými třídami a **logických pravidel**;
- ✓ pomocí **vzdálenosti od reprezentativních obrazů** (etalonů) klasifikačních tříd;
- ✓ pomocí **ztotožnění s etalony**;

# OBECNÉ SCHÉMA ZPRACOVÁNÍ DAT



# OBECNÉ SCHÉMA ZPRACOVÁNÍ DAT

## UČENÍ

### ☑ **učení klasifikátoru**

→ nastavení klasifikačních kritérií;

☐ s učitelem

- dokonalým
- nedokonalým

☐ bez učitele – typicky shlukování

### ☑ **výběr prvků popisu dat**

→ stanovení reprezentativních charakteristických rysů zpracovávaného dat;

# TYPY KLASIFIKÁTORŮ

Základní členění vychází z reprezentace vstupních dat

- ☑ **příznakové** – každý vstupní data jsou vyjádřena vektorem hodnot (příznaků);
  - paralelní (např. Bayesův klasifikátor, ...)
  - sekvenční (např. klasifikační stromy, ...)
- ☑ **strukturální (syntaktické)** – vstupní data jsou popsána relačními strukturami;
- ☑ **kombinované** – jednotlivá primitiva jsou doplněna příznakovým popisem



# TYPY KLASIFIKÁTORŮ

## Deterministický klasifikátor

- každá deterministická klasifikace musí být jednoznačná a úplná, tzn., že každý obraz (předmět, jev) musí patřit do nějaké třídy a nemůže být současně ve dvou či více třídách.

## Pravděpodobnostní klasifikátor

- pravděpodobnostní klasifikátor stanoví pravděpodobnost zařazení obrazů do daných klasifikačních tříd

# TYPY KLASIFIKÁTORŮ

Na základě typů klasifikačních a učících algoritmů:

- ✓ parametrické;
- ✓ neparametrické

# KLASIFIKACE x PREDIKCE

**predikce** (z lat. *prae-*, před, a *dicere*, říkat) zjevně nese časové hledisko, když jej používáme ve významu předpověď či prognózu, jako soud o tom, co se stane nebo nestane v budoucnosti. V tomto významu je používán např. v analýze či zpracování časových řad.

(prediction x forecasting)

# KLASIFIKACE x PREDIKCE

pojem **klasifikace** je používán, použije-li se klasifikačního algoritmu pro známá data. Pokud jsou data nová, pro která apriori neznáme klasifikační třídu, pak hovoříme o predikci klasifikační třídy.

<http://www.kdnuggets.com/faq/classification-vs-prediction.html> (23.8.2010)

# KLASIFIKACE x PREDIKCE

pojem **klasifikace** používáme, pokud vybíráme identifikátor klasifikační třídy z určitého diskrétního konečného počtu možných identifikátorů. Pokud určujeme (predikujeme) spojitou hodnotu, např. pomocí regrese, pak hovoříme o predikci, i když tento pojem nemá časovou dimenzi.

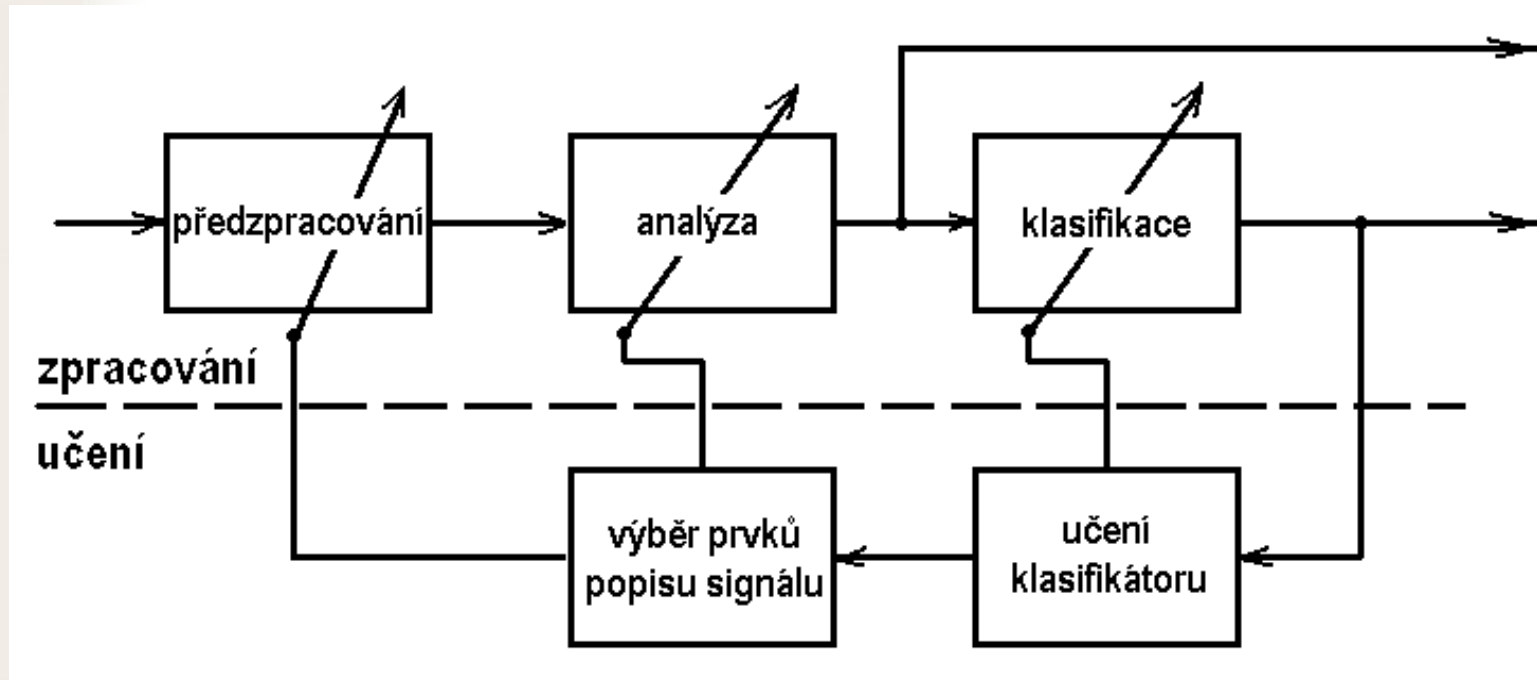
Han, J., Kamber, M.: Data Mining Concepts and Techniques. The Morgan Kaufmann Series in Data Management Systems. 2<sup>nd</sup> edition, Elsevier; Amsterdam(2005), 800 s.

# DISKRIMINAČNÍ ANALÝZA

týká se obecně vztahu mezi kategoriální proměnnou a množinou vzájemně vázaných příznakových proměnných.

Konkrétně, předpokládejme že existuje konečný počet, řekněme  $R$ , různých a priori známých populací, kategorií, tříd nebo skupin, které označujeme  $\omega_r$ ,  $r=1, \dots, R$  a úkolem diskriminační analýzy je nalézt vztah, na základě kterého pro daný vektor příznaků popisujících konkrétní objekt tomuto vektoru přiřadíme hodnotu  $\omega_r$ .

# OBECNÉ SCHÉMA ZPRACOVÁNÍ DAT



# ZÁVĚREM SHRNU TÍ

- ✓ co je to klasifikace?
- ✓ klasifikace vs. predikce vs. diskriminační analýza
- ✓ základní principy klasifikace
- ✓ parametrická vs. neparametrická klasifikace



# Příprava nových učebních materiálů pro obor Matematická biologie

je podporována projektem ESF  
č. CZ.1.07/2.2.00/07.0318

## „VÍCEOBOROVÁ INOVACE STUDIA MATEMATICKÉ BIOLOGIE“



INVESTICE DO ROZVOJE VZDĚLÁVÁNÍ