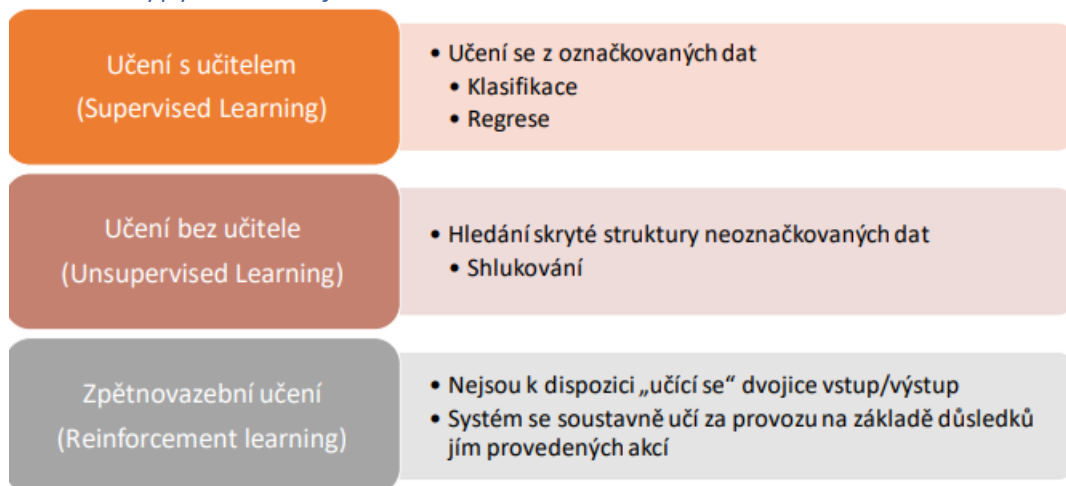


## 9. Základní typy úloh strojového učení a rozdělení dat. Základy klasifikace, vyhodnocování výsledků, matice záměn.

### Základní typy úloh strojového učení a rozdělení dat



### Klasifikace

- Klasifikace znamená doslova třídění respektive zařazování do různých tříd
- Klasifikační algoritmy se učí ve fázi trénování na základě dat, u kterých je známa jejich příslušnost k daným předem definovaným třídám
- Při nasazení naučeného klasifikátoru se pak určuje pro každý klasifikovaný objekt, do které třídy spadá

#### Celá řada typů klasifikátorů, nejčastější:

- Rozhodovací stromy (nebudeme se dále zabývat, příklad viz. dále)
- Vzdálenostní
- Pravděpodobnostní
- Logistická regrese
- Support Vector Machine (SVM)
- **Neuronová síť typu vícevrstvý perceptron (Multi Layer Perceptron – MLP)**

### Regrese

- Umožňuje odhadovat hodnotu jedné náhodné veličiny (závislé proměnné) na základě znalosti jiných veličin (nezávislých proměnných)
- Regresní algoritmy se opět učí ve fázi trénování na základě dat, u kterých známe hodnoty závislé veličiny i nezávislých veličin
- Na základě těchto dat lze vytvořit příslušný regresní model
- Ve fázi nasazení systému je tento model využit pro predikci hodnot závislé proměnné pro jiné (další) hodnoty nezávislých proměnných

- Lze měřit přesnost predikce – například jako kvadratickou odchylku predikované a skutečné hodnot

## Rozdělení dat

### Trénovací data

- Slouží ke stanovení parametrů modelu

### Ověřovací (validační) data

- Nalezení optimálních parametrů modelu (počet parametrů, struktura modelu, ...)
- Lze vyčlenit z trénovacích dat a provést tzv. cross-validaci

### Testovací data

- Nesmí být použita ke stejnému účelu jako trénovací nebo validační data
- Umožňují pouze vyhodnotit přetrénování nebo generalizační schopnosti modelu

Trénovací data

Ověřovací  
(validační)  
data

Testovací data

## Základy klasifikace

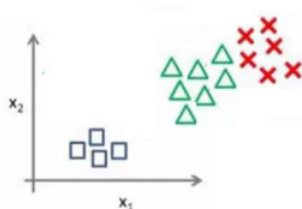
### Učení s učitelem

=Během učení se klasifikátoru jsou v příznakovém prostoru vymezeny **diskriminační funkce** (definují oblasti, v kterých všechny body (objekty) náleží k dané třídě)

- Klasifikovaný objekt lze reprezentovat jako vektor příznaků s rozměrem  $D$
- Např. jak bude vypadat vektor příznaků pro černobílý obrázek a jaká bude
- Každý objekt respektive jeho vektorová reprezentace určuje v daném  $D$ -dimenzionálním prostoru příznaků jeden bod
- Prostory o velkém počtu dimenzí nelze vykreslovat na 2D ploše  $\Lambda$
- Počet dimenzí si proto pro další výklad omezíme na dvě až tři

## Prostor příznaků – 2D příklad

- Pro objekty, popsané jen 2 příznaky,  $x_1$  a  $x_2$ , je prostor příznaků 2D rovina
- Jednotlivé třídy se pak znázorňují různými symboly a/nebo barvami
- Objekty patřící do jednotlivých tříd mají podobné hodnoty příznaků
  - Vytváří v prostoru příznaků shluky (jsou blízko sebe)



- Klasifikátor osob na třídy **Dítě**, **Muž**, **Žena**

- $x_1$ ...výška,  $x_2$ ...hmotnost

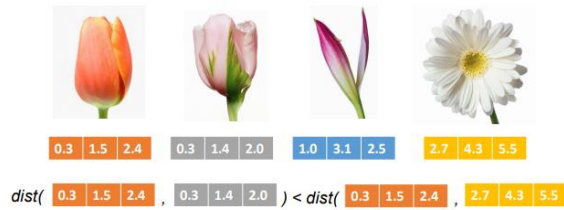
- Čtverečky...třída Děti

- Křížky... třída Muži

- Trojúhelníky...třída Ženy

## Vzdálenostní metody klasifikace

Jsou-li objekty  $x_1$  a  $x_2$  podobné, je i jejich reprezentace v příznakovém prostoru blízká a jejich vzdálenost malá.



## NN a KNN

### Metoda nejbližšího souseda

#### Učení

- Uložení všech  $N$  trénovacích dat
- Známe hodnoty dat a příslušnost každého vzorku k nějaké třídě

#### Klasifikace do celkem $K$ tříd $c_1 \dots c_K$

- Pro klasifikovaný příznakový vektor  $x$  je nalezen vzorek  $(x_i, c_i)$ , který je mu nejbližší, a na základě jeho příslušnosti k jedné ze tříd je rozhodnuto o klasifikaci

$$x \in c_k \leftrightarrow dist(x, x_i) < dist(x, x_j), x_i \in c_k \quad 1 \leq j \leq N$$



Ethalon = reprezentant třídy

## Vzdálenostní funkce

Euklidovská vzdálenost (Euclidean dist.)

$$d(x, z) = \sqrt{\sum_{i=1}^p (x_i - z_i)^2}$$

Vzdálenost v městských blocích (Manhattan dist.)

$$d(x, z) = \sum_{i=1}^p |x_i - z_i|$$

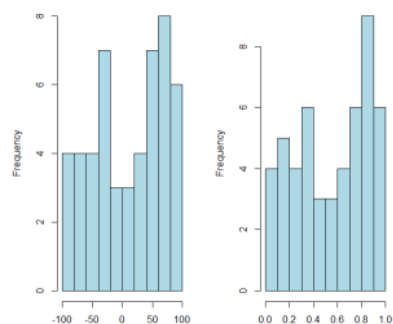
## Normalizace příznaků

### Normalizace příznaků – metoda Max-Min

Normalizovaná hodnota  $z_i = \frac{x_i - \min(x)}{\max(x) - \min(x)}$

Hodnoty normalizovaných dat leží v intervalu  $<0,1>$

Nejsou zachovány záporné hodnoty a hodnoty větší než 1



## Normalizace příznaků – standardizace (z-skóre)

Normalizovaná hodnota  $z_i = \frac{x_i - \text{mean}(x)}{\sigma(x)}$

Jsou zachovány kladné i záporné hodnoty

Normalizovaná data mají nulovou střední hodnotu a jednotkový rozptyl

- Používá se pro trénování neuronových sítí

### Klasifikace v prostoru obecně

- Obecně každý klasifikátor vymezuje v prostoru příznaků rozdělovací podplochy, které určují hranice jednotlivých tříd • Obtížně znázornitelné na 2D ploše pro více dimenzí než 2 nebo 3
- Tyto hranice mezi jednotlivými třídami představují tzv. diskriminační funkce
- Podle typu diskriminační funkce lze klasifikátory dělit
- Základní dělení je na **lineární** (A) a **nelineární** (B)

### Matice záměn

		Skutečná třída	
		+	-
Predikovaná třída	+	True Positive (TP)	False Positive (FP)
	-	False Negative (FN)	True Negative (TN)

Obrázek 4.2: Matice záměn binární klasifikace

### Vyhodnocování výsledků

Lze vyhodnotit následujícími parametry:

- Precision
- Recall
- F1 skóre
- mAP (AP)
- ROC křivka