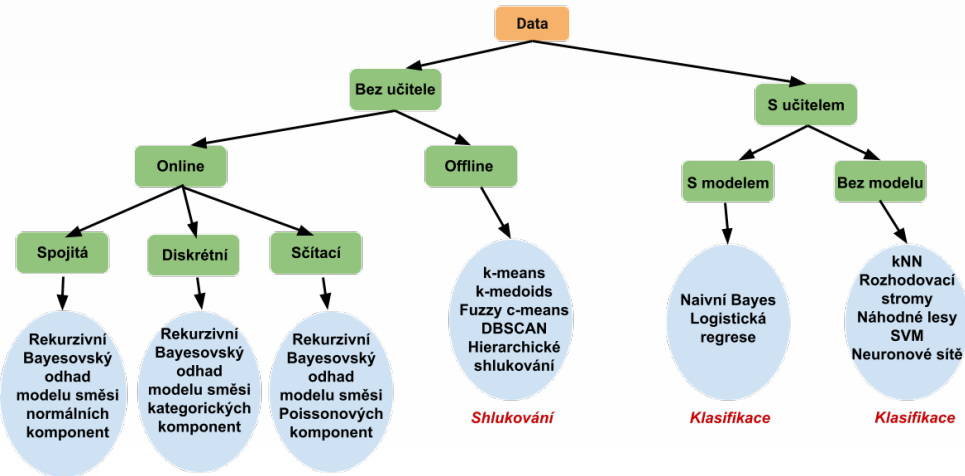


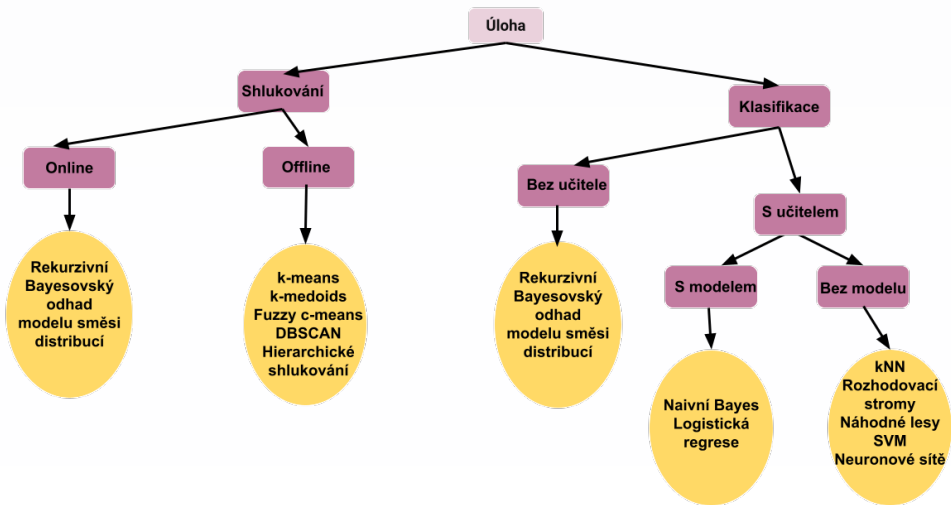
# Přednáška 9 – Shrnutí a opakování algoritmů

## Jak zvolit vhodný algoritmus podle dat?



*Shlukování a klasifikace*

# Jak zvolit vhodný algoritmus podle úlohy?



# Příklady – který algoritmus je vhodný?

- 1 Měříme:** tep a vzdálenost na hodinkách.  
**Co chceme:** detekovat v reálném case, zda se jedná o klidový režim, chůzi nebo běh.
- 2 Měříme:** hladinu hluku ve sluchátkách (pravé, levé).  
**Co chceme:** detekovat v reálném case, zda hladina zvuku je bezpečná.
- 3 Máme data:** věk a cena automobilových náhradních dílů zařazených do kategorií “nové originální”, “neoriginální”, “použité originální”.  
**Co chceme:** zařadit do kategorií další náhradní díly
- 4 Měříme:** počty návštěvníků webu za minutu a jejich lokalitu.  
**Co chceme:** určit, zda aktuální provoz webu je normální nebo se jedná o zvýšenou aktivitu.
- 5 Měříme:** hodnoty z NO<sub>x</sub> senzoru před a po NO<sub>x</sub> katalyzátoru v autě.  
**Co chceme:** za provozu vozidla vyhodnotit, zda se jedná o správnou funkci katalyzátoru nebo poruchu.



- komponenty – modely jednotlivých shluků  $f_j(y_t|\Theta), j = 1, \dots, n_c$
- ukazovátka – diskretní náhodná veličina  $c_t \in \{1, \dots, n_c\}$ ,  
ukazuje aktivní komponentu

## Základní princip – obecně:

- Inicializace komponent

### V cyklu:

- Měříme data
- Vzdálenosti od jednotlivých komponent – **proximity**
- Normalizace proximit – **váhy komponent**  
(pravděpodobnost, že je komponenta aktivní)
- Maximální váha – **bodový odhad ukazovátka**  
(klasifikace)
- Update statistik **s váženými daty**
- Přepočítání bodových odhadů podle typu komponent
- Jdeme na krok 1

# Otázky 30 – 32, 35 – 40: shlukování a klasifikace online s modelem směsi normálních komponent

- **spojitá** data, každá komponenta – **normální** model
- generování:  $c(t) = \text{sum}(\text{rand}(1,1) > \text{cumsum}(a1)) + 1$ ;  
 $y(:,t) = \text{th}(:,c(t)) + \text{sqrt}(r) * \text{randn}(2,1)$ ;
- vzdálenost – rozptyly      • inicializace – histogram, x-y graf

Pro čas  $t = 0$

- 1 Počet komponent, počáteční bodové odhady  $(\hat{\theta}_0)_j$ ,  $(\hat{r}_0)_j$ , statistiky

Pro čas  $t = 1, 2, \dots$ , pro každou komponentu

- 1 Měříme nová data  $y_t$
- 2 Proximity  $m_j$  – dosadíme do komponent  $(\hat{\theta}_{t-1})_j$ ,  $(\hat{r}_{t-1})_j$  a  $y_t$
- 3 Váhy komponent  $w_{j;t} = \frac{m_j}{\sum_{i=1}^{n_c} m_i}$
- 4 Bodový odhad ukazovátka:  $\hat{c}_t = \arg \max_j w_t$
- 5 Update statistik s váženými daty:

$$(V_t)_j = (V_{t-1})_j + w_{j;t} \begin{bmatrix} y_t \\ 1 \end{bmatrix} [y_t' \ 1], \quad (\kappa_t)_j = (\kappa_{t-1})_j + w_{j;t}$$

- 6 Přepočítání bodových odhadů  $(\hat{\theta}_t)_j$ ,  $(\hat{r}_t)_j$

- **diskrétní** data, každá komponenta – **kategorický** model
- generování: 
$$c(t) = \text{sum}(\text{rand}(1, 1) > \text{cumsum}(a_l)) + 1;$$
$$pp = \text{cumsum}(\text{th}\{c(t)\});$$
$$y(t) = \text{sum}(\text{rand}(1, 1) > pp) + 1;$$

## Algoritmus

Pro čas  $t = 0$

- 1 Počet komponent, počáteční statistiky *expertně*, bodové odhady  $(\hat{\Theta}_0)_j$

Pro čas  $t = 1, 2, \dots$ , pro každou komponentu

- 1 Měříme nová data  $y_t$
- 2 Proximity  $m_j$  – dosadíme do komponent  $(\hat{\Theta}_{t-1})_j$  a  $y_t$
- 3 Váhy komponent  $w_{j;t} = \frac{m_j}{\sum_{i=1}^{n_c} m_i}$
- 4 Bodový odhad ukazovátka:  $\hat{c}_t = \arg \max_j w_t$
- 5 Update statistik s **váženými daty**:  $(\nu_{i;t})_j = (\nu_{i;t-1})_j + w_{j;t} \delta(i; y_t)$
- 6 Přepočítání bodových odhadů  $(\hat{\Theta}_{i;t})_j = \frac{(\nu_{i;t})_j}{\sum_{k=1}^n (\nu_{k;t})_j}$

## Otázky 45 – 51: shlukování a klasifikace online s modelem směsi Poissonových komponent

- **sčítací data**, každá komponenta – **Poissonův** model
- generování: 

```
c(t)=sum(rand(1,1)>cumsum(a1))+1;
for i=1:ny y(i,t)=poissrnd(La{i}(c(t))) end
```
- vzdálenost –  $\lambda$       • inicializace – expertně, histogram, x-y graf

### Algoritmus

Pro čas  $t = 0$

- 1 Počet komponent, počáteční bodové odhady  $(\hat{\lambda}_0)_j$ , statistiky

Pro čas  $t = 1, 2, \dots$ , pro každou komponentu

- 1 Měříme nová data  $y_t$
- 2 Proximity  $m_j$  – dosadíme do komponent  $(\hat{\lambda}_{t-1})_j$  a  $y_t$
- 3 Váhy komponent  $w_{j;t} = \frac{m_j}{\sum_{i=1}^{n_c} m_i}$
- 4 Bodový odhad ukazovátka:  $\hat{c}_t = \arg \max_j w_t$
- 5 Update statistik s **váženými daty**:  
 $(S_t)_j = (S_{t-1})_j + w_{j;t} y_t$ ,  $(\kappa_t)_j = (\kappa_{t-1})_j + w_{j;t}$ ,
- 6 Přepočítání bodových odhadů  $(\hat{\lambda}_t)_j = \frac{(S_t)_j}{(\kappa_t)_j}$



# Otázky 52 – 55: Poissonova regrese (predikce sčítacích dat)

$y_t$  – Poissonovo rozdělení

$$f(y_t | \lambda) = \exp\{-\lambda\} \frac{\lambda^{y_t}}{y_t!}$$

Poissonova regrese

$$\ln(\underbrace{\bar{y}_t}_{\lambda}) = x_{1;t}\theta_1 + x_{2;t}\theta_2 + \dots + k$$

$\lambda$  = průměr = rozptyl

$x_t = [x_{1;t} \ x_{2;t} \ \dots \ x_{n;t}]$  – vysvětlující data

Linearizace + nejmenší čtverce

metoda 1

$$\underbrace{\begin{bmatrix} \ln(y_1) \\ \ln(y_2) \\ \dots \\ \ln(y_{nd}) \end{bmatrix}}_Y = \underbrace{\begin{bmatrix} x_{1;1} & x_{2;1} & \dots & 1 \\ x_{1;2} & x_{2;2} & \dots & 1 \\ \dots & \dots & \dots & \dots \\ x_{1;nd} & x_{2;nd} & \dots & 1 \end{bmatrix}}_X \underbrace{\begin{bmatrix} \theta_1 \\ \theta_2 \\ \dots \\ k \end{bmatrix}}_\theta$$

Odhad offline:

trénovací data

$$\hat{\theta} = (X'X)^{-1}X'Y$$

Predikce online:

testovací data

$$\hat{Y} = \exp\{X\hat{\theta}\}$$

Metoda maximální věrohodnosti

metoda 2

$$\ln L_t(\theta) = \ln \prod_{\tau=1}^t \text{model}(\text{trénovací data}_\tau), \quad \hat{\theta} = \arg \max_{\theta} \ln L_t(\theta)$$

- s učitelem – známé shluky (ukazovátka)

## Naivní Bayesův klasifikátor:

- vektor veličin – jednotlivé veličiny – **podmíněná nezávislost**
- model:

$$f([y_{1;t}, y_{2;t}] | \Theta, c_t) = \prod_{i=1}^{N=2} f(y_{i;t} | \Theta_i, c_t)$$

### • Postup:

1. **Trénování modelů** – trénovací data (s učitelem – měřené ukazovátka)
  - odhad parametrů modelů (filtrace každé veličiny + offline odhad)
2. **Testování (klasifikace)** – testovací data (bez ukazovátka)

• odhad ukazovátka z odhadnutých modelů

$$f(c_t | y_t) = \underbrace{\text{normovaný histogram ukazovátka}}_{f(c_t)} \times \underbrace{\text{součin modelů}}_{\prod_{i=1}^{N=2} f(y_{i;t} | \hat{\Theta}_i, c_t)}$$

- model:

- pro  $c_t \in \{0, 1\}$ :

$$\frac{\exp\{c_t z_t\}}{1 + \exp\{z_t\}}, \quad z_t = y_{1;t}\theta_1 + y_{2;t}\theta_2 + \dots + k$$

- pro  $c_t \in \{0, 1, \dots, N\}$ :

$$\ln \frac{p_1}{p_0} = z_{t(1)}, \ln \frac{p_2}{p_0} = z_{t(2)} \dots, \ln \frac{p_N}{p_0} = z_{t(N)}$$

## Postup:

### 1. **Trénování modelů** – trénovací data (s učitelem – měřené ukazovátko)

- odhad regresních koeficientů – maximální věrohodnost

$$\ln L_t(\theta) = \ln \prod_{\tau=1}^t \text{model}(\text{trénovací data}_\tau), \quad \hat{\theta} = \arg \max_{\theta} \ln L_t(\theta)$$

### 2. **Testování (klasifikace)** – testovací data (bez ukazovátko)

- odhad ukazovátko

# Otázky 62 – 63: vyhodnocení přesnosti klasifikace

Příklad: binární klasifikace (pozitivní, negativní)

- Chybová matice (confusion matrix)

|                    |                    |                    |
|--------------------|--------------------|--------------------|
|                    | Skutečně pozitivní | Skutečně negativní |
| Pozitivní predikce | TP                 | FP                 |
| Negativní predikce | FN                 | TN                 |

- Přesnost (accuracy)

Procento správných predikcí ze všech dat

$$Acc = \frac{TP + TN}{TP + FP + FN + TN}$$

Pro nevyvážená data:

- Preciznost (precision)

Procento správných pozitivních predikcí mezi pozitivními predikcemi

$$Pre = \frac{TP}{TP + FP}$$

- Výtěžnost (recall)

Procento správných pozitivních predikcí mezi skutečně pozitivními

$$Rec = \frac{TP}{TP + FN}$$

- Skóre F1  $\in (0, 1)$

$$F1 = 2 \frac{Pre \times Rec}{Pre + Rec}$$

## Otázky 64 – 69: shlukovací metody offline bez učitele (k-means, k-medoids, fuzzy, dbscan, hierarchické)

- bez učitele – hledáme shluky (nemáme ukazovátka)
- na základě vzdálenosti mezi datovými body

Euklidovská: 
$$m_E(x, y) = \sqrt{(x_1 - y_1)^2 + \dots + (x_n - y_n)^2}$$

### k-means

- vzdálenost (bod, středy)
- nejbližší střed – shluk
- průměr bodů ve shluku = **nový střed**

### k-medoids

- vzdálenost (bod, medoidy)
- nejbližší medoid – shluk
- vzdálenost uvnitř shluků
- medoid za ne-medoid

### Fuzzy c-means

- středy shluků
- matice příslušnosti

### DBSCAN

- $\epsilon$ -okolí, minimální počet bodů
- shluk: vnitřní bod + všechny z něj dosažitelné
- hledáme hustě propojený shluk

### Hierarchické shlukování

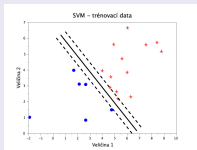
- aglomerativní – “zdola nahoru”, divizivní – “shora dolů”
- propojení – min/max/průměrná vzdálenost

- s učitelem – máme shluky (máme ukazovátka)

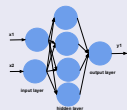
## kNN

- vzdálenost (nový bod, body)
- $k$  nejbližších bodů
- nový bod  $\in$  do shluku většiny nejbližších

## Podpůrné vektorové stroje (SVM)



## Neuronové sítě



- aktivační funkce
- vícevrstvé perceptrony
- zpětné šíření chyby

## Rozhodovací stromy



- IG  $\rightarrow$  1 – lepší strom

## Náhodné lesy

- Bootstrap Aggregation
- 10 – 100 náhodných vzorků
- les = 10 – 100 stromů
- každý strom v lese se učí
- kombinace výsledků stromů

## Rekurzivní Bayesovské odhadování modelu směsi distribucí:

- I. Nagy, E. Suzdaleva. Algorithms and Programs of Dynamic Mixture Estimation. Unified Approach to Different Types of Components, Springer, 2017.
- Kárný, M., Böhm, J., Guy, T. V., Jirsa, L., Nagy, I., Nedoma, P., Tesař, L., 2006. Optimized Bayesian Dynamic Advising: Theory and Algorithms. Springer-Verlag, London.

## Shlukování a klasifikace offline:

- D. T. Larose. Discovering Knowledge in Data. An Introduction to Data Mining. Willey, 2005.
- J. Han, M. Kamber, J. Pei. Data Mining: Concepts and Techniques, 3rd Edition. Morgan Kaufmann, 2011.
- M. J. Zaki, W. Meira Jr. Data Mining and Analysis: Fundamental Concepts and Algorithms. Cambridge University Press, 2014.
- G. James, D. Witten, T. Hastie, R. Tibshirani, J. Taylor. An Introduction to Statistical Learning with Applications in Python. Springer Texts in Statistics. Springer Cham, 2023.