



# ZZN – Řešení

3. 1. 2018

A. Kvapilová (xkvapi12), T. Coufal (xcoufa09)

# Dolovací úloha

## Míra znečištění v závislosti na teplotě a vlhkosti

Úloha předpokládá závislost míry znečištění na okolních podmínkách a zkoumá trendy, které změny počasí provází. Na základě výsledovaných závislostí lze určit míru rizik spojených především s denní dobou, teplotou a vlhkostí. Míra znečištění je určena na základě evropských norem pro kvalitu ovzduší.

## Rozbor úlohy

Určení míry znečištění je definováno pomocí tzv. AQI indexu. Ten je v různých zemích počítán jinak. V Evropské unii je tento index zjišťován z šesti různých faktorů<sup>1</sup>:

- 8h průměr koncentrace CO
- 8h průměr koncentrace O<sub>3</sub>
- aktuální hodinová koncentrace NO<sub>2</sub>
- aktuální denní koncentrace PT<sub>2,5</sub>
- aktuální denní koncentrace PT<sub>10</sub>
- aktuální denní koncentrace SO<sub>2</sub>

V daném datasetu byly ovšem k dispozici pouze první tři zmiňované. Naopak jsou k dispozici hodnoty<sup>2</sup> koncentrací těžkých kovů, benzenu a NMHC a dalších, pro které ovšem existují pouze roční bezpečné limity – neexistuje specifikace pro určení kvality ovzduší na základě jejich momentálních hodnot. Z toho důvodu nebyly při výpočtu uvažovány.

Takto získaný AQI index definuje kvalitu ovzduší a míru znečištění jako kategorie: **Good, Moderate, Unhealthy, Very Unhealthy, Hazardous, Very Hazardous**.<sup>3</sup>

V řešené úloze jsme tuto kategorii znečištění predikovali pouze na základě teploty, absolutní vlhkosti, relativní vlhkosti, data a času, tedy pro takový systém, který by nemusel nutně měřit koncentrace látek v ovzduší.

## Schéma řešení

1. Filtrace a předzpracování dat odstraněním irelevantních atributů a chybějících hodnot
2. Výpočet 8h průměrů pro CO a O<sub>3</sub> (*Moving Average* z rozšíření *Value Series*)
3. Klasifikace do AQI kategorií podle evropských norem (*Generate Attributes, Map, Set Role*)
4. Výpočet dílčích klasifikátorů, pro lepší predikci (*Generate Attributes*)
5. Klasifikace znečištění pomocí Naive Bayes (*Validation, Select Attributes, Naive Bayes, Apply Model, Performance*)

Pro lepší výsledky predikce jsme vedle zřejmých charakteristik jako teplota a vlhkost přidali klasifikátory pro:

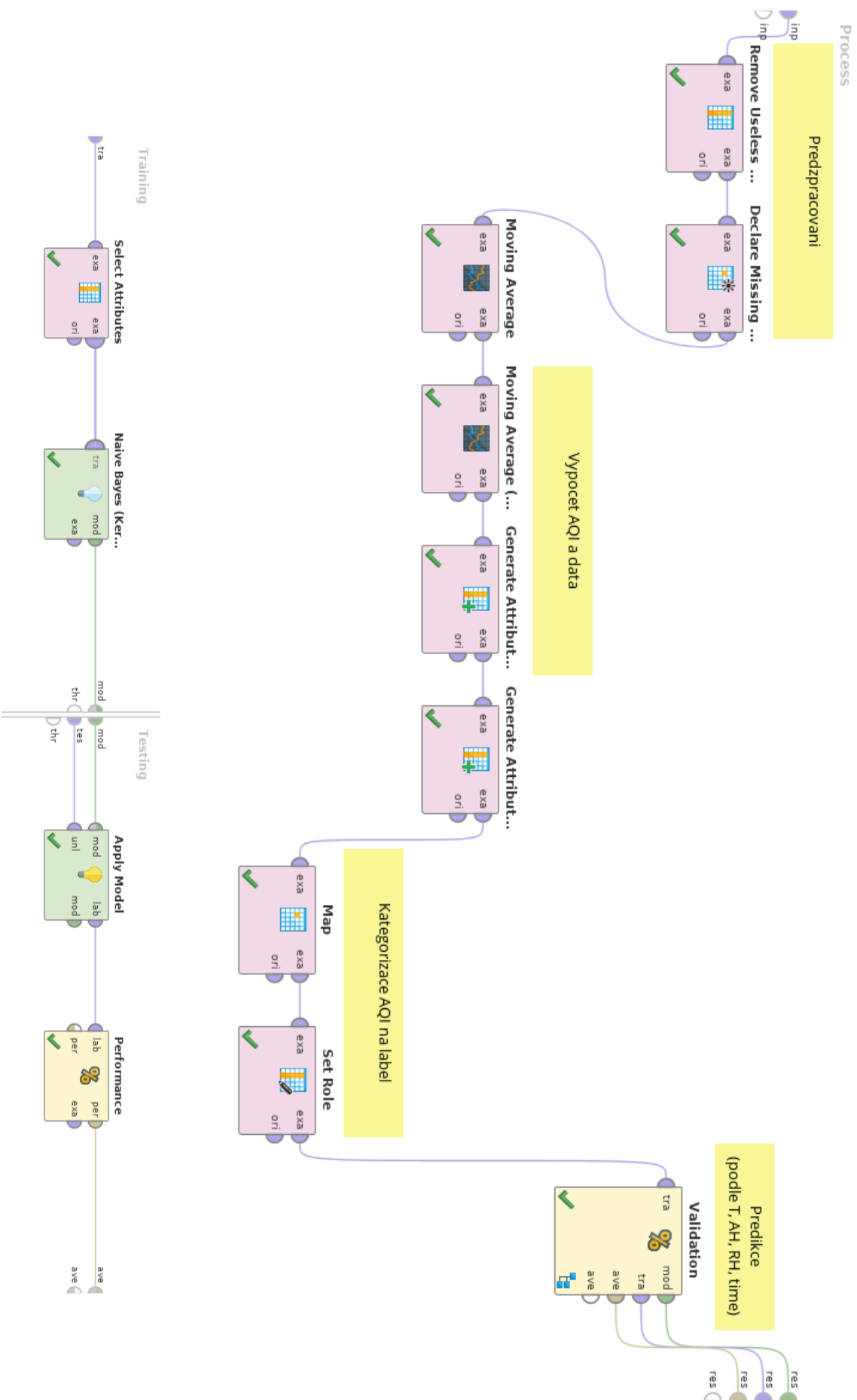
- Den v týdnu
- Denní doba (noc, siesta, práce)
- Roční období
- Pracovní dny (resp. dny s charakteristicky nižší produkcí škodlivin: Sobota, Neděle a Pondělí)
- Hodina měření

Pro naučení prediktivního modelu jsme použili Naive Bayes klasifikaci, která je vhodná pro menší soubory dat a navíc dobře klasifikuje numerické hodnoty.

<sup>1</sup> [https://www.airqualitynow.eu/about\\_indices\\_definition.php](https://www.airqualitynow.eu/about_indices_definition.php)

<sup>2</sup> <https://archive.ics.uci.edu/ml/datasets/Air+quality>

<sup>3</sup> [http://www.haze.gov.sg/docs/default-source/computation-of-the-pollutant-standards-index-\(psi\).pdf](http://www.haze.gov.sg/docs/default-source/computation-of-the-pollutant-standards-index-(psi).pdf)



## Závěr

Výsledný model predikuje míru znečištění s 60% úspěšností. To je způsobeno především nízkou kvalitou vstupních dat:

- chybějící faktory znečištění
- vychýlením dat směrem k vysokému znečištění
- mnoho chybějících hodnot
- malý soubor dat (1 rok pozorování)
- Nedostatek meteorologických údajů (např. tlak)

I přesto byla klasifikace relativně úspěšná a pokud nebyl testovací vzorek zařazen do správné kategorie, byl modelem přiřazen do vedlejší (tzn. Pokud měl být vzorek Hazardous byl chybně zařazen jako Very Hazardous), nikdy však nedošlo k jeho zařazení do opačné části spektra znečištění.

Predikovaná klasifikace	Správná klasifikace			Class precision
	Very Unhealthy	Hazardous	Very Hazardous	
Very Unhealthy	600	173	253	58.48%
Hazardous	124	92	67	32.51%
Very Hazardous	278	184	892	65.88%
Class recall	59.88%	20.49%	73.60%	