

# **Tematické modelování publikační činnosti České akademie věd a umění v letech 1890–1910**

Diplomová práce

# Tematické modelování publikační činnosti České akademie věd a umění v letech 1890–1910

## Cíl práce:

Na základě strojového analýzy tištěných publikací zjistit, jaká témata byla předmětem bádání České akademie císaře Františka Josefa pro slovesnost, vědy a umění v prvních dvaceti letech jejího fungování.

## Obsah:

### Teoretická část

1. Česká akademie věd a umění
2. Digitalizace fondů
3. Makroanalýza digitálních dat

### Výzkumná část

4. Modelování témat
5. Diskuze

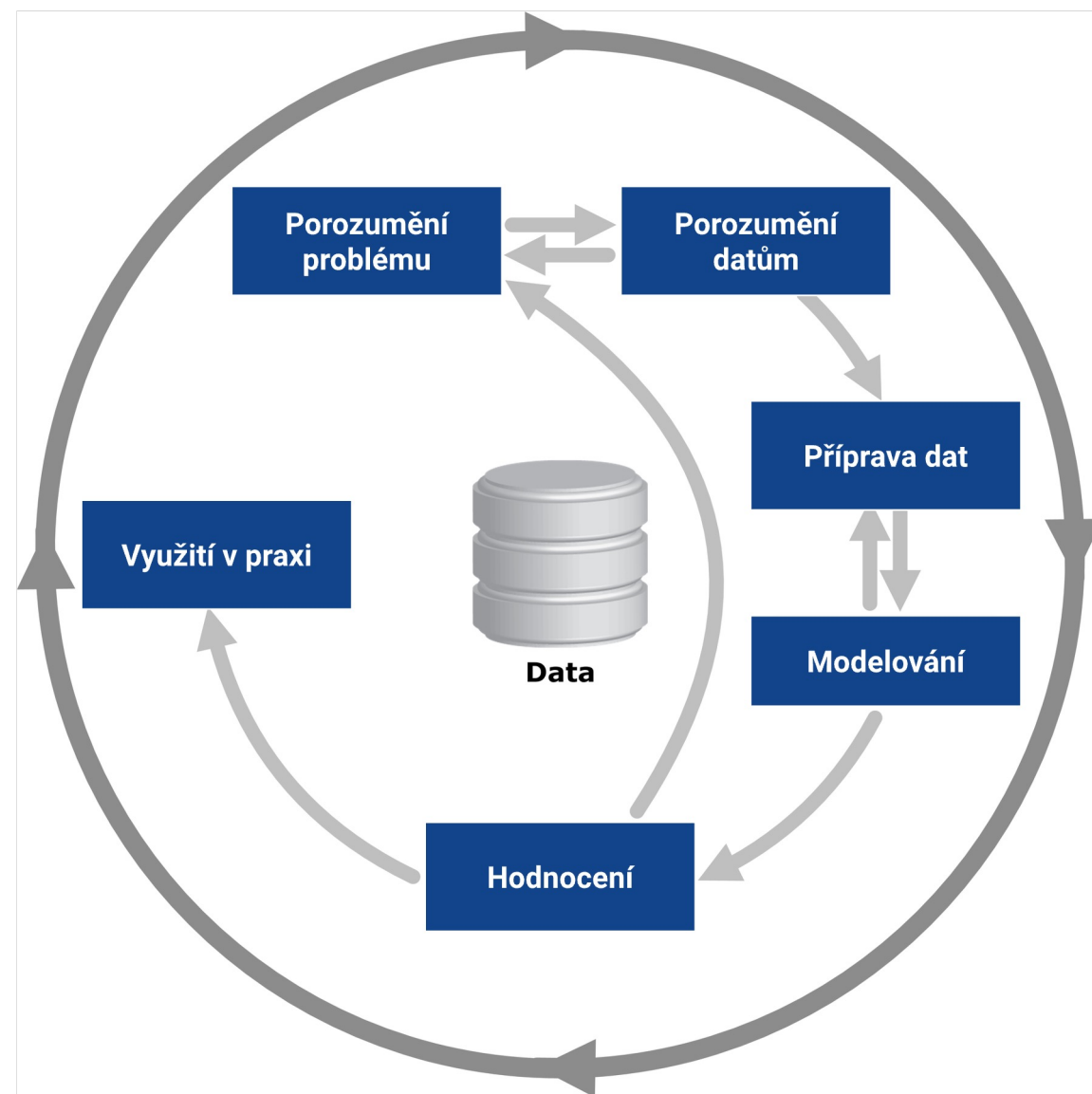
# Česká akademie věd a umění

- Česká akademie císaře Františka Josefa pro vědy, slovesnost a umění vznikla 23. ledna 1890
- Organizována byla do 4 tříd, tato skutečnost měla dopad na publikační činnost:
  - společné publikace (*Věstník, Almanach*),
  - každá třída vydávala **Rozpravy**,
  - monografie,
  - podpora vydávání časopisů.
- **Témata (dosavadní stav poznání):**
  - I. třída: historie, právní problematika a filozofie,
  - II. třída: matematika, fyzika, chemii, biologie (bryologie), lékařství, technika,
  - III. filologie, čeština, česká literatura, lexikografie, dialektologie, orientální jazyky,...

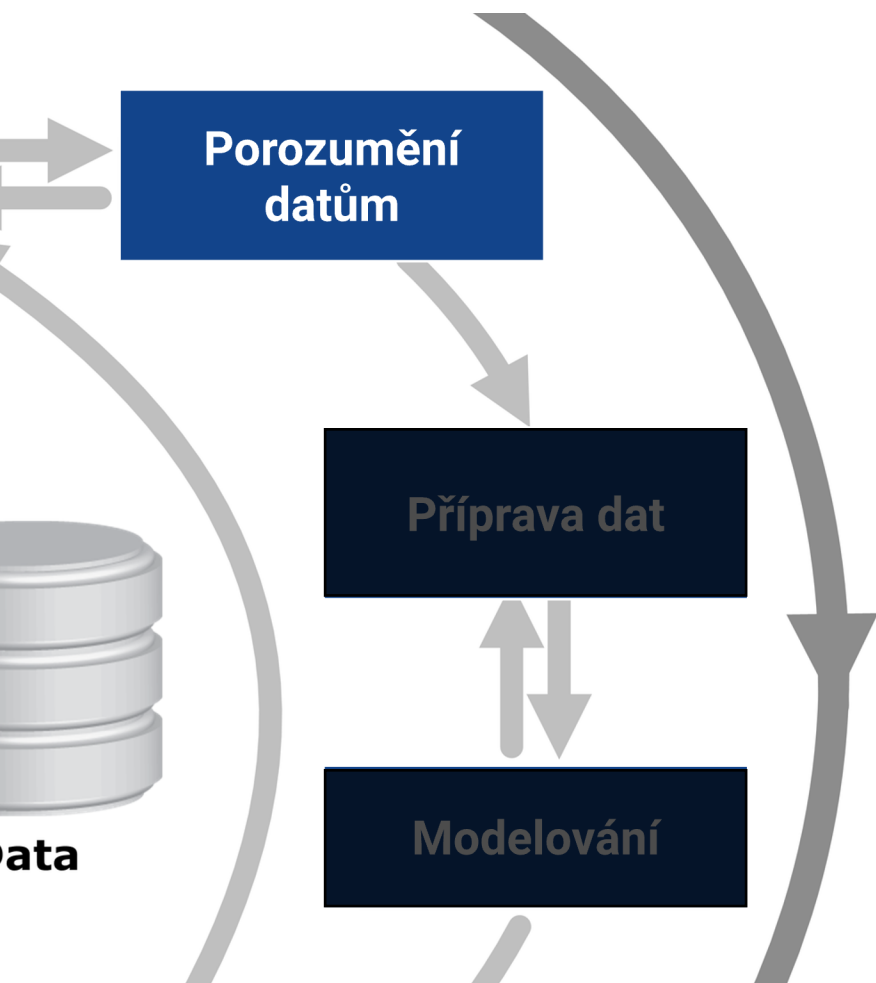
# Digitalizace fondů a makroanalýza digitálních dat

- Reformátování (digitalizace) *Rozprav* představuje nový zdroj pro analýzu.
- Distant reading / makroanalýza:
  - Na texty v digitálních knihovnách lze nahlížet jako na data, ve kterých lze identifikovat vzorce a souvislosti.
- Tematické modelování
- **Latentní Dirichletova alokace**

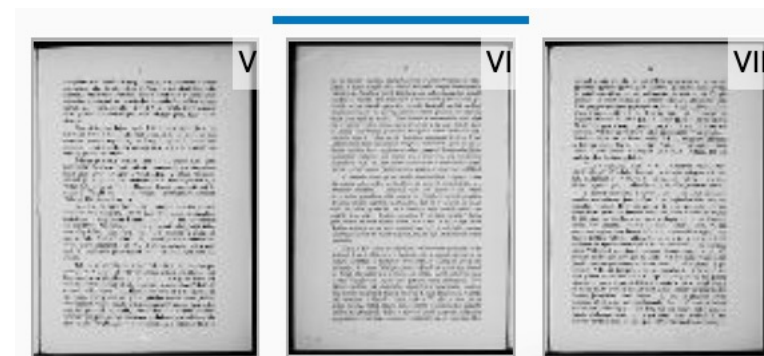
# Výzkumná část



# Porozumění datům



Digitální knihovna Akademie věd ČR

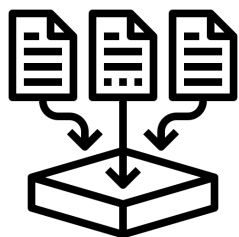
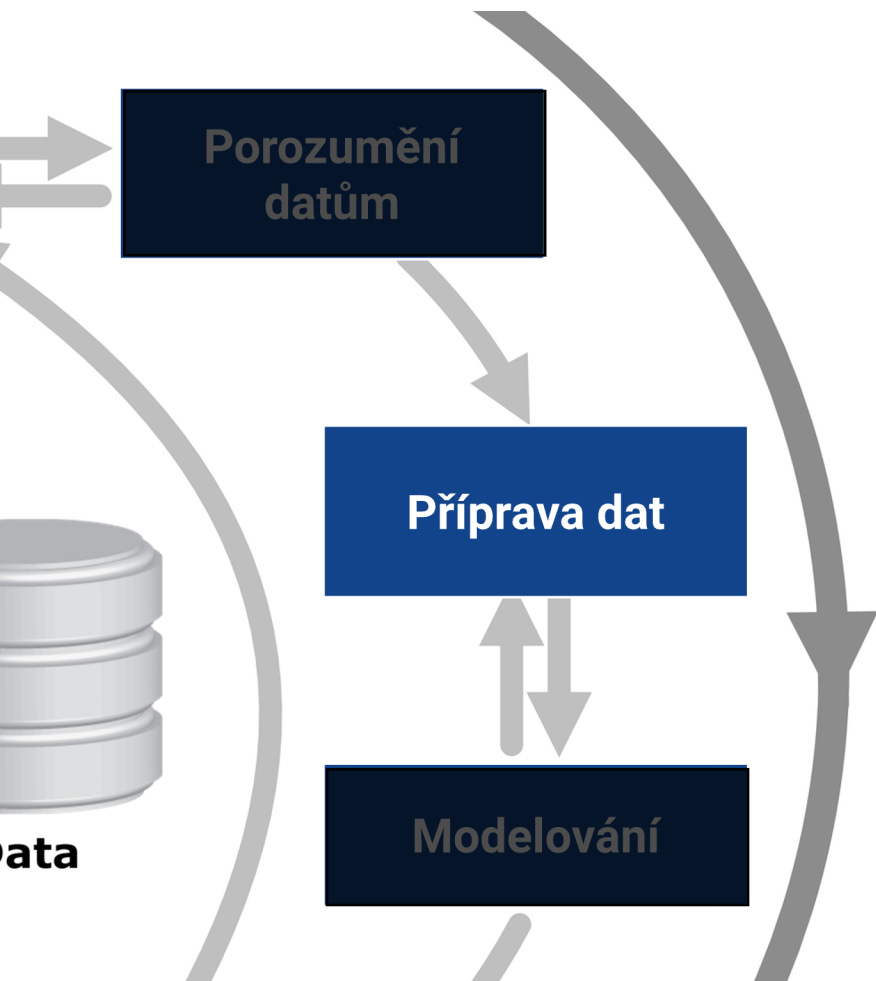


[ocr](#) | [stránka](#) | číslo periodika | ročník periodika | noviny [URL](#)

## PŘEDMLUVA.

Roku 1878 vydal jsem »Vergilstudien nebst einer Collatio  
V mnohých recensích vysloveno přání a očekávání, že poda  
Toto pokračování podávám nyní, kdy zřízením České Akader  
V práci této podány jsou k velikému počtu míst Vergilio  
•'-1 Seznal jsem ovšem také tu i tam, ale celkem ve vel  
l\*

# Příprava dat



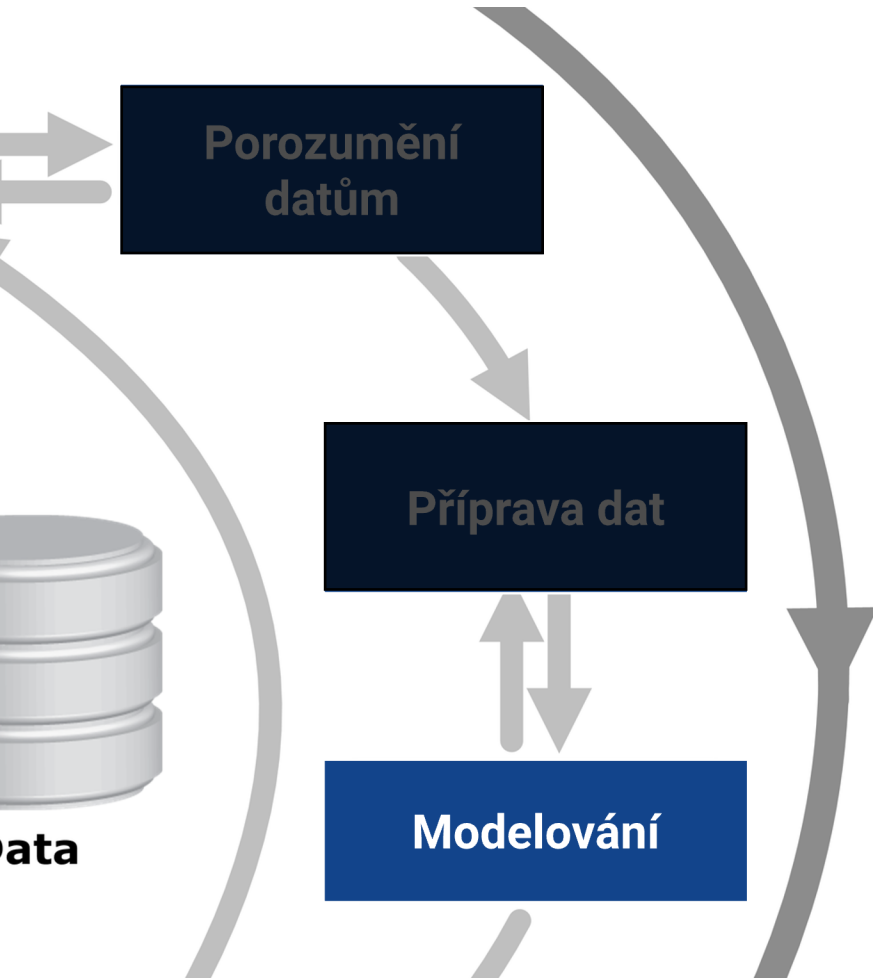
**.txt soubory:**  
887 svazků,  
31 147 stran

## Transformace, úpravy, obohacení:

- odstraněna slova obsahující dva nebo méně znaků,
- obohacení pomocí UDPipe 2 (tokenizace, lemmatizace a identifikaci slovního druhu),
- vyfiltrována jen podstatná jména.

**Pro každý svazek *Rozprav* vznikl jeden .txt soubor obsahující pouze podstatná jména v základním tvaru**

# Modelování



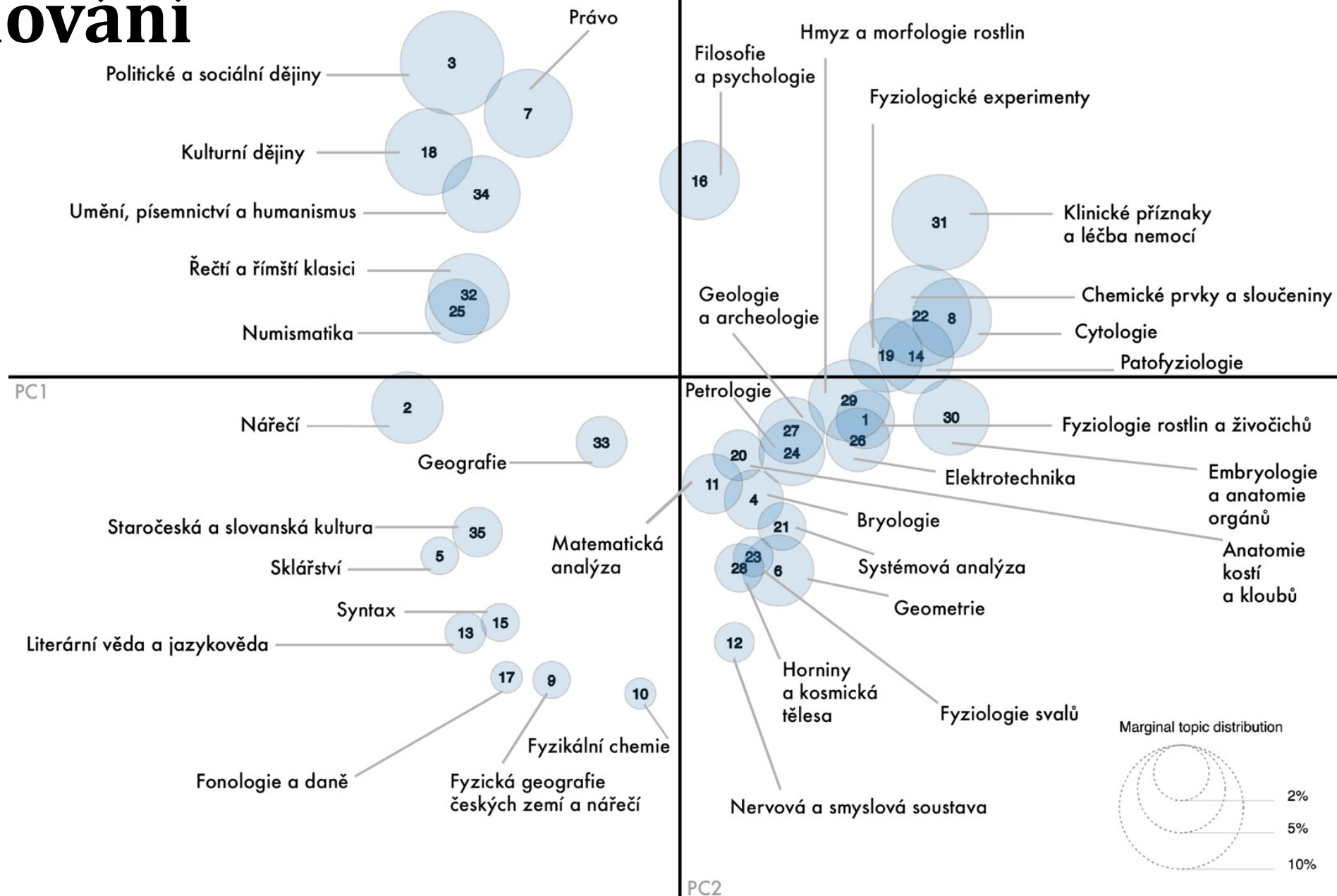
- Vstupem pro model LDA je **matice dokument–slovo**
- K samotnému modelování témat byla využita knihovna `topicmodels`
  - Vyžaduje zvolit počet témat, která mají být v korpusu identifikována.
  - Na základě kvalifikovaného odhadu bylo zvoleno 35 témat.
  - **Cílem modelování je identifikovat témata s alespoň o úroveň jemnější granularitou, než doposud.**



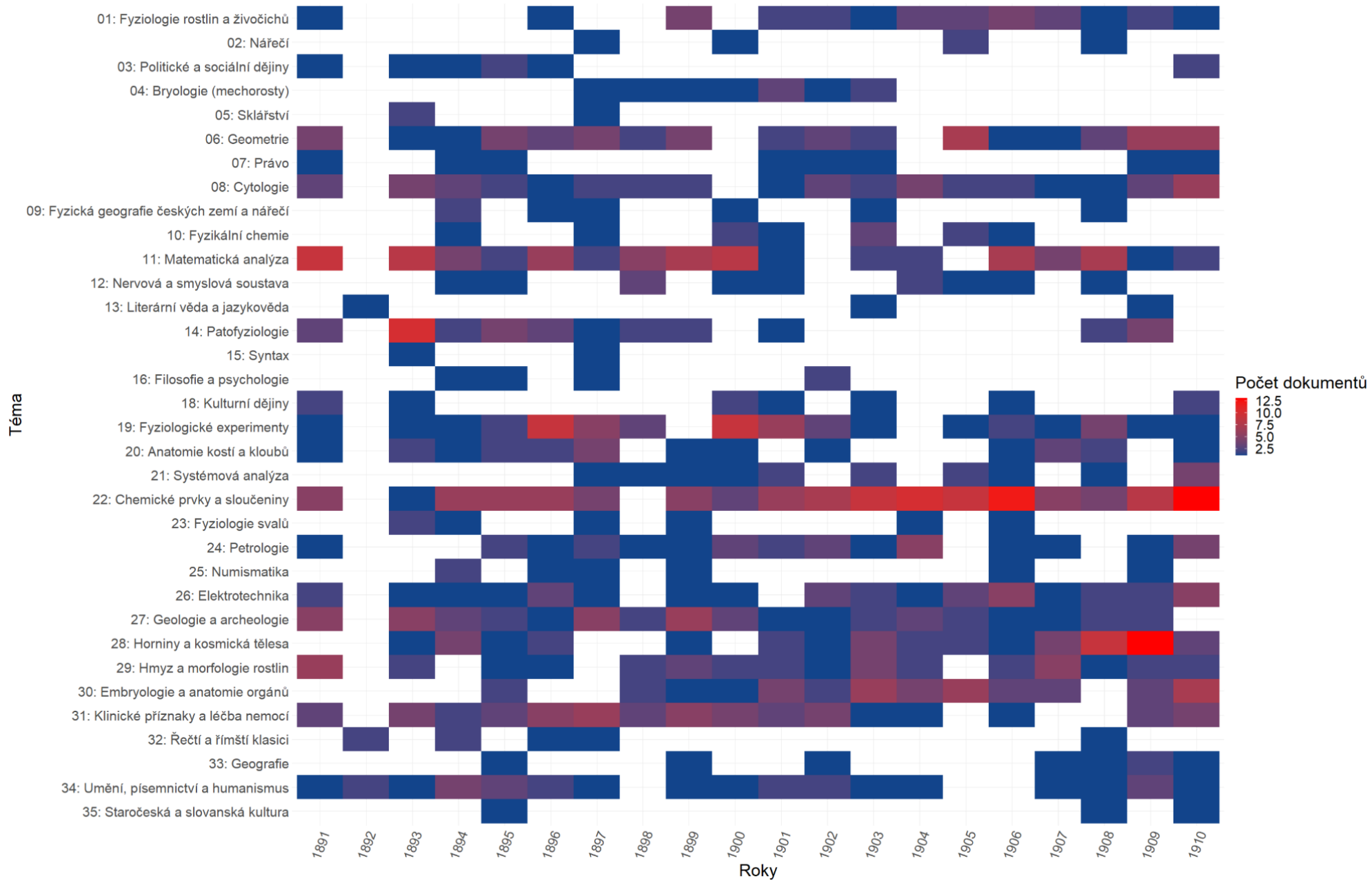
# Modelování

- Výstupem modelu jsou dvě matice:
  - **slovo–téma**, ve které je u každého slova z korpusu uvedeno, s jakou pravděpodobností „patří“ k jakému tématu.
  - **dokument–téma** popisující pravděpodobnost příslušnosti dokumentu ke konkrétnímu tématu.
- **Model rozpoznává témata v korpusu označuje pouze číslem**, pro pojmenování byly využity tři techniky:
  - nejfrekventovanější slova, která se v každém tématu vyskytují (17 případů),
  - analýza začátků dokumentů patřících k danému tématu s nejvyšší pravděpodobností,
  - analýza zdrojových digitalizátů.

# Modelování



# Modelování



# Výsledky

- Provedená analýza podává přehled **konkrétních témat, kterým se Česká akademie ve sledovaném období věnovala**,
    - jejich sémantické podobnosti,
    - počtu svazků zařazených do každého z témat
    - a proměn počtu vydaných svazků v jednotlivých letech.
  - Rozpoznaná témata byla namapována na kategorie vycházející z vědních sekcí AV ČR a také vědních oborů a oborových skupin FORD.
  - Práce **představuje podrobný pohled na publikační činnost České akademie jako celek** a umožňuje rozšířit dosavadní stav poznání oblastí zájmů této instituce.
-

# Závěr

- Součástí diplomové práce jsou 3 přílohy:
  - **Příloha č. 1:** Seznam publikací vydaných Českou akademií věd císaře Františka Josefa pro vědy, slovesnost a umění z období let 1890–1918.
  - **Příloha č. 2:** Postup stažení a přípravy dat a modelování témat
    - Data a scripty byly zveřejněny také na platformě GitHub: [https://github.com/kerschfilip/tematicke\\_modelovani\\_cavu](https://github.com/kerschfilip/tematicke_modelovani_cavu)
    - Veškerá data jsou dostupná v repozitáři Zenodo: <https://doi.org/10.5281/zenodo.10395970>
  - **Příloha č. 3:** Seznam odkazů na digitalizované svazky časopisu Rozpravy rozřazený podle rozpoznaných témat

# **Tematické modelování publikační činnosti České akademie věd a umění v letech 1890–1910**

Diplomová práce