

MINISTERSTVO PRO MÍSTNÍ ROZVOJ  
Národní orgán pro koordinaci

# Data, vizualizace a infografiky pro podporu core aktivit TA ČR

Systémy, principy a výstupy

Martin Víta, 3. 10. 2018



EVROPSKÁ UNIE  
Fond soudržnosti  
Operační program Technická pomoc



MINISTERSTVO  
PRO MÍSTNÍ  
ROZVOJ ČR

**T A**

**Č R**

**Výzkum užitečný pro společnost**

**T A**  
**Č R**



Evropská unie  
Evropský sociální fond  
Operační program Zaměstnanost

# **Data, vizualizace a infografiky pro podporu core aktivit TA ČR**

**Systémy, principy a výstupy**

3. 10. 2018

**Technologická agentura České republiky**

**Martin Víta**

Metodik QC/QA - Projekt ProEval

- **Monitoring projektů / programů**
- **Evaluaace programů**
- Zkvalitňování veřejných soutěží
- Ad hoc analýzy pro potřeby TA ČR a dalších institucí
- ...

*Existence relevantních datových zdrojů + trend  
evidence based policy: ideální kombinace*

# T A Data ve výzkumu, vývoji a inovacích v ČR Č R

- **Klíčová role IS VaVaI** (dnes na rvvi.cz)
  - problém: omezená funkcionality týkající se vyhledávání/filtrování
- **Interní data TA ČR**
- Oficiální statistika ČSÚ a EUROSTAT
- Ostatní zdroje (ČR a EU)

**Obecný problém integrace datových zdrojů**

*Systém STARFOS jakožto nástroj na integraci dat,  
interaktivní filtrování a vyhledávání*

T A  
Č R

STARFOS

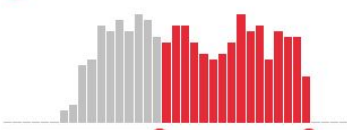
strojové učení

CS  
EN

English

Roky řešení

☒ Průběh ☐ Začátek ☐ Konec



2002 - 2017 249

Subjekt

Finance

☒ Celkové náklady ☐ Podpora

☐ do 500 tis. (53)


☐ 500 - 999 tis. (62)

☒ 1 000 - 2 499 tis. (86)

☐ 2 500 - 4 999 tis. (106)

☐ 5 000 - 9 999 tis. (139)

Nalezeno 86 projektů za 0,306s




GAP103/10/1875: Učení se z teorií  
Současné algoritmy **strojového učení** se snaží budovat obecné modely (teorie ... **strojového učení**, v němž se agenti učí nejen z pozorování, ale především ... ). Za tímto účelem využijeme současných technologií...  
1 989 tis. Kč | 1 989 tis. Kč | GA ČR | 2010–2012

GA201/09/1665: Překonání propasti mezi systémovou biologii a strojovým učním  
Překonání propasti mezi systémovou biologii a strojovým učním vytvořením algoritmů schopných těžit znalosti z procesních sítí, jako jsou metabolické, genově-regulační či buněčně-signální dráhy, které představují...  
1 983 tis. Kč | 1 983 tis. Kč | GA ČR | 2009–2011

GA201/08/0509: Integrace strojového učení a splňování omezujících podmínek  
Navrhujeme integrovat pokročilé techniky dvou vyspělých výzkumných oborů, a to relačního strojového učení (RSU) na jedné straně a splňování omezujících podmínek (SOP) na straně druhé ... , a to relačního **strojového**...  
1 938 tis. Kč | 1 938 tis. Kč | GA ČR | 2008–2010

GA201/05/0557: Aproximace a učení funkcí více proměnných pomocí neuronových sítí a jadrových metod  
Cílem projektu je přispět k interakci mezi klasickými a nově vytvářenými matematickými pojmy a rozvíje ... matematickými pojmy a rozvíjejícími se obory neuropočítání a **strojového učení**. Aby mohly ... aflexibilní...  
1 260 tis. Kč | 1 260 tis. Kč | GA ČR | 2005–2007

ID15113: Aplikace umělé inteligence v astronomii



- **Datová základna** (na počátku data rvvi.cz)
- **Fasetová klasifikace** (a vyhledávání, resp. filtrování)
  - *system třídění umožňující zařazení objektu do více tříd (výzkumná organizace / příjemce podpory od daného poskytovatele / instituce z daného kraje / ...)*
  - *obvyklý např. ve velkých e-shopech*
- **Search engine** (solr)

- Systém **integrující různorodé relevantní datové zdroje**
- Komfortní aplikace na **vyhledávání a filtrování** entit různých typů ve VaVal (instituce, výsledky, projekty ...)
- Nástroj na tvorbu **interaktivních reportů** – inspirace D3 (data driven documents)



**T A**  
**Č R**

## **Část II – vizualizace**

**[www.tacr.cz](http://www.tacr.cz)**

- Motivace – jednoduché otázky
  - Jak získat přehled o zaměření projektů daného programu v “kompaktní formě”?
  - Které instituce jsou z pohledu daného programu “významné?” - viz další část přednášky
- Nutnost přejít od reportů k vizualizacím
- *Ukázka (shlukování na množině projektů: separátní PDF)*

# Principy práce s entitami reprezentovanými textově

- Předpokládáme, že entity máme popsány pomocí **textových dokumentů** (např. instituce seznamy klíčových slov z jejich projektů)
- **Vektorová reprezentace** textů
- Výpočet jejich **podobnosti** (vzdálenosti ve vektorovém prostoru)
- Analýza **grafu podobnosti** (např. detekce komunit/clusterů)

- Využití popsaného přístupu **v jiném kontextu**
- Prakticky motivovaná otázka: “Jaké jsou typické tématické okruhy dotazů uchazečů?”
- *Ukázky vizualizací (separátní soubory):*
  - *shluky/komunity otázek*
  - *wordcloudy nejčastějších slov/slovních spojení*

T A  
Č R

## Část III – analýza grafových dat

[www.tacr.cz](http://www.tacr.cz)

- Otázka souvýskytů určitých entit, např. klasifikačních tříd oborů (CEP)
  - *hlavní a vedlejší obor projektu*
- Výstup – opět graf podobnosti
- *Ukázka (separátní PDF)*

- Motivační otázky:
  - Které instituce jsou z pohledu daného programu “významné”?
  - Vznikají v rámci programu kliky či komunity spolupracujících institucí (bez apriorní negativní konotace)
- *Ukázka grafu (separátní PDF)*

# T A Modelování významnosti entit – centrality

## Č R

- Hlavní myšlenky jednotlivých typů centralit:
  - degree centrality
  - betweenness centrality
  - closeness centrality
  - triangle
  - **eigenvector centrality**



T A  
Č R

# Institute s nejvyšší eigenvector centrality v první VS programu ALFA

Institute	Eigenvector centrality
České vysoké učení technické v Praze	10.0
Vysoké učení technické v Brně	0.53
Vysoká škola báňská - Technická univerzita Ostrava	0.32
Centrum dopravního výzkumu, v.v.i.	0.26
Ústav teorie informace a automatizace AV ČR, v. v. i.	0.21
Západočeská univerzita v Plzni	0.20
EVEKTOR, spol. s r.o.	0.19
HVM PLASMA, spol. s r.o.	0.19
VARS BRNO a.s.	0.18

**T A**  
**Č R**

## **Část IV – TA ČR a infografiky**

**[www.tacr.cz](http://www.tacr.cz)**

- Vizualizace jakožto vhodně zvolená grafická reprezentace konkrétních dat
- Infografika jakožto souhrn klíčových informací o dané problematice grafickou a interaktivní formou
- *Procházka infografikami <http://visual.tacr.cz>*

**T A**

**Č R**

# **Část V – datové zdroje, TA ČR a otevřená data**

**[www.tacr.cz](http://www.tacr.cz)**

- **INKA** – mapování inovačních kapacit
  - <https://inkaviz.tacr.cz> (prohlížeč dat)
- **TA ČR v číslech** - statické prostředí, které nás navedlo k visualu
  - ...další zdroje na <https://tacrcz.cz/test-programy/1158-datove-zdroje-vav.html>

- Strojová čitelnost (požadavky na formát)
- Nevylučující a jasně daná licence
- Katalogizace

# T A První kroky TA ČR v oblasti otevřených dat Č R

- Příprava publikování prvních datových sad: nejdůležitější: **stav projektů TA ČR** (napříč veřejnými soutěžemi)
  - Otevřenost: stupeň 3
  - Katalogizace – “katalog as a service”
- Později přibudou podkladové datové sady k infografikám, data ze systému INKAviz

- Využívání otevřených dat z dalších zdrojů (např. ČSSZ, ČSÚ – doplňování mozaiky dat např. o krajích)
- Využívání otevřených dat ostatních poskytovatelů podpory ve VaVal (“které instituce neúspěšně podávající projekty v loňském roce se hlásí nyní k nám do veřejné soutěže...?”)

*Nutnost vývoje aplikací nad otevřenými daty a spolupráce s dalšími institucemi...*



**T A  
Č R**

**Děkuji za pozornost**

**[www.tacr.cz](http://www.tacr.cz)**