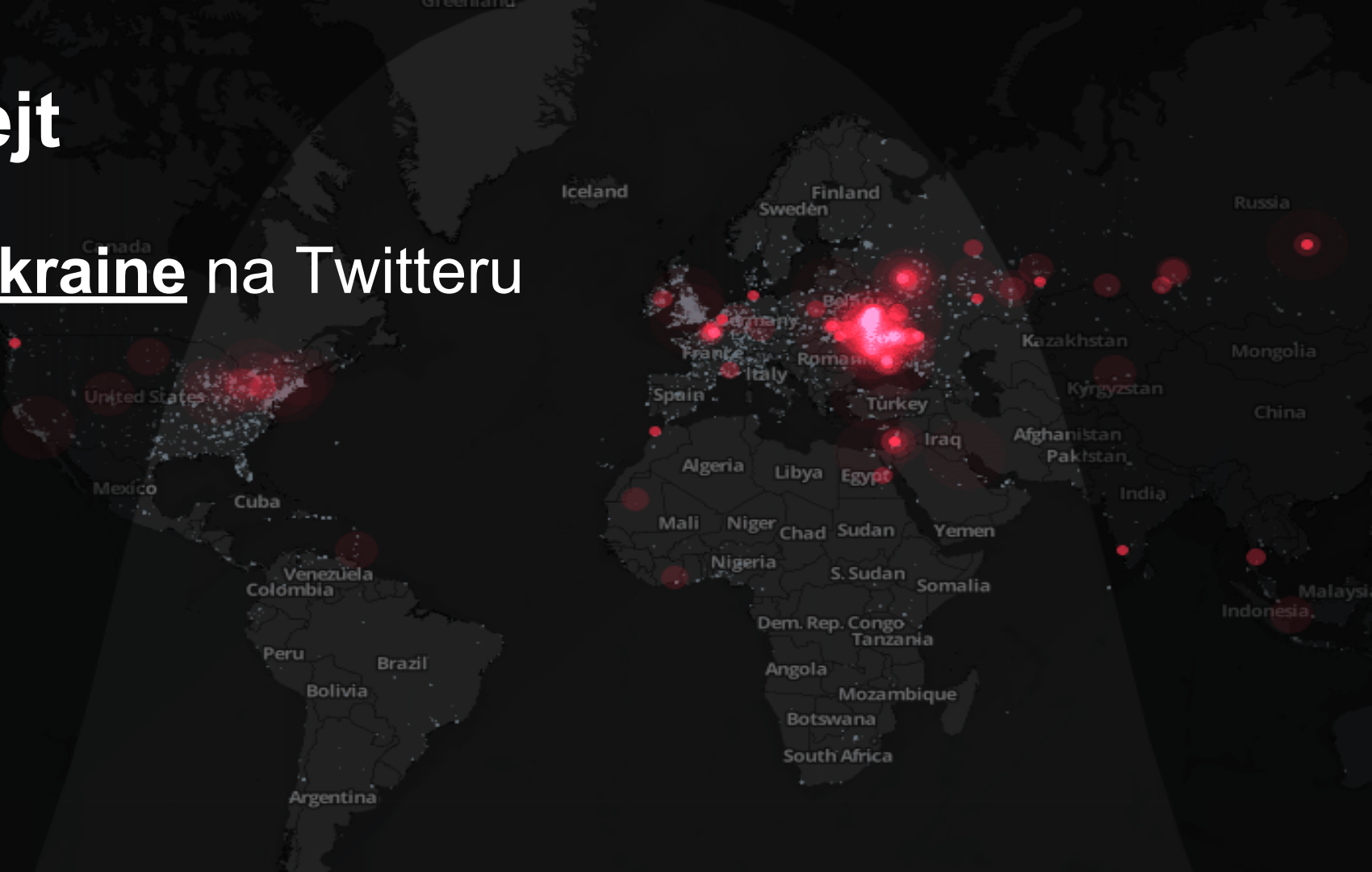


Zdroje dat

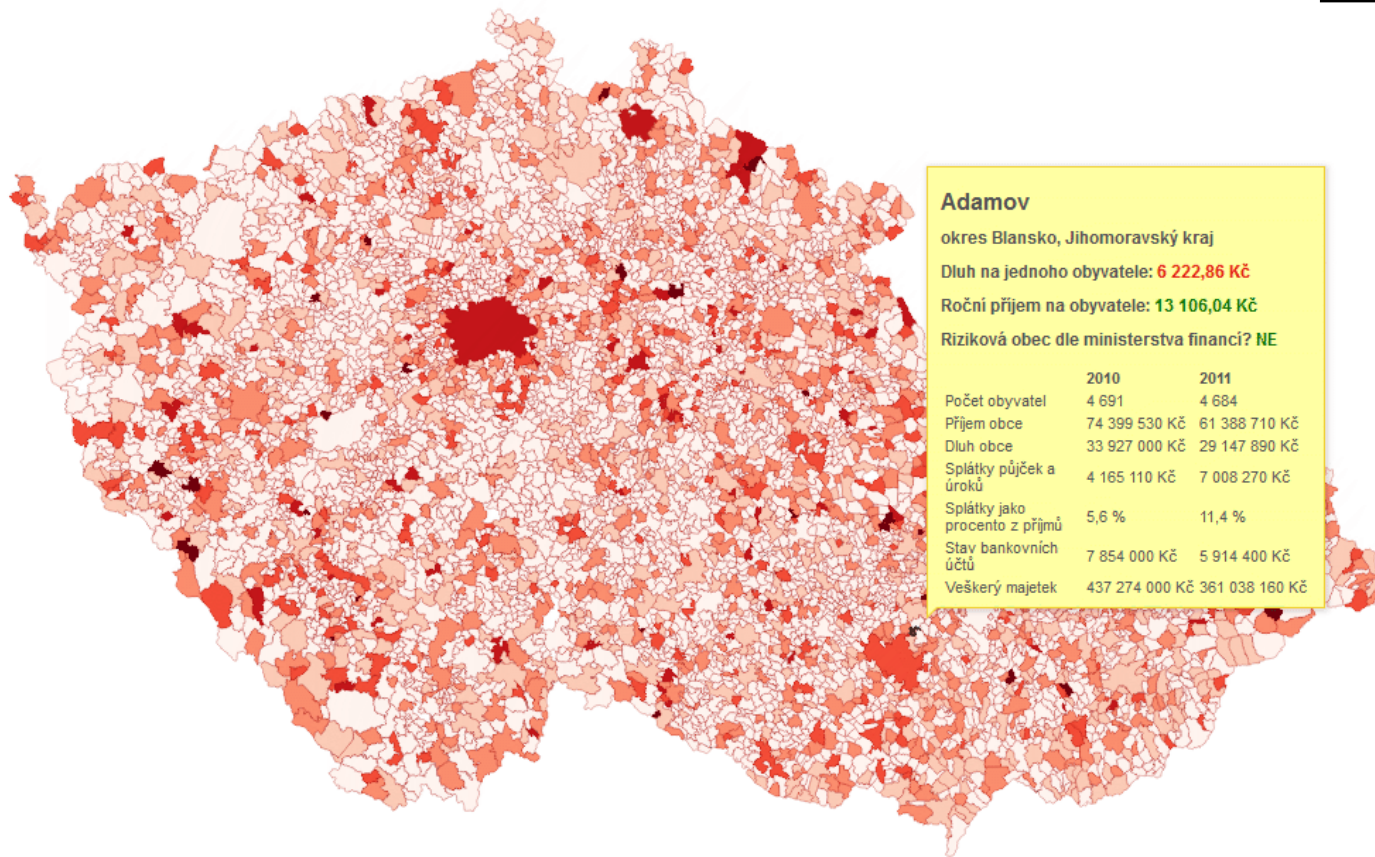
Hejt

#Ukraine na Twitteru



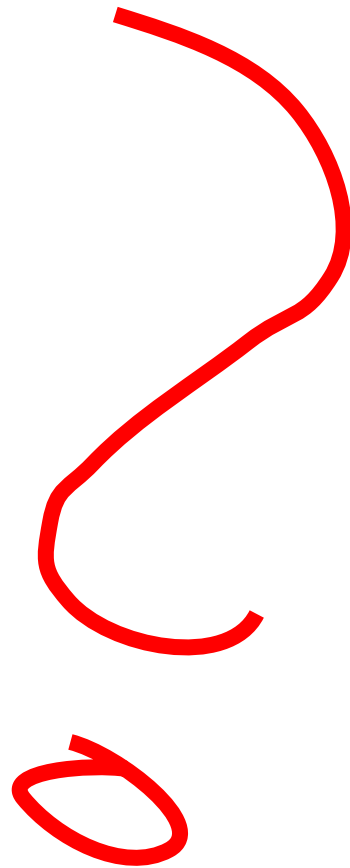
Hejt

- Mapa dluhů
- Mapa Ukrajiny
- Bezpečnost na D1



Plán

- Zdroje a formáty dat - přednáška
- Čištění dat - seminář
- Scrapování webu - seminář



Zdroje dat

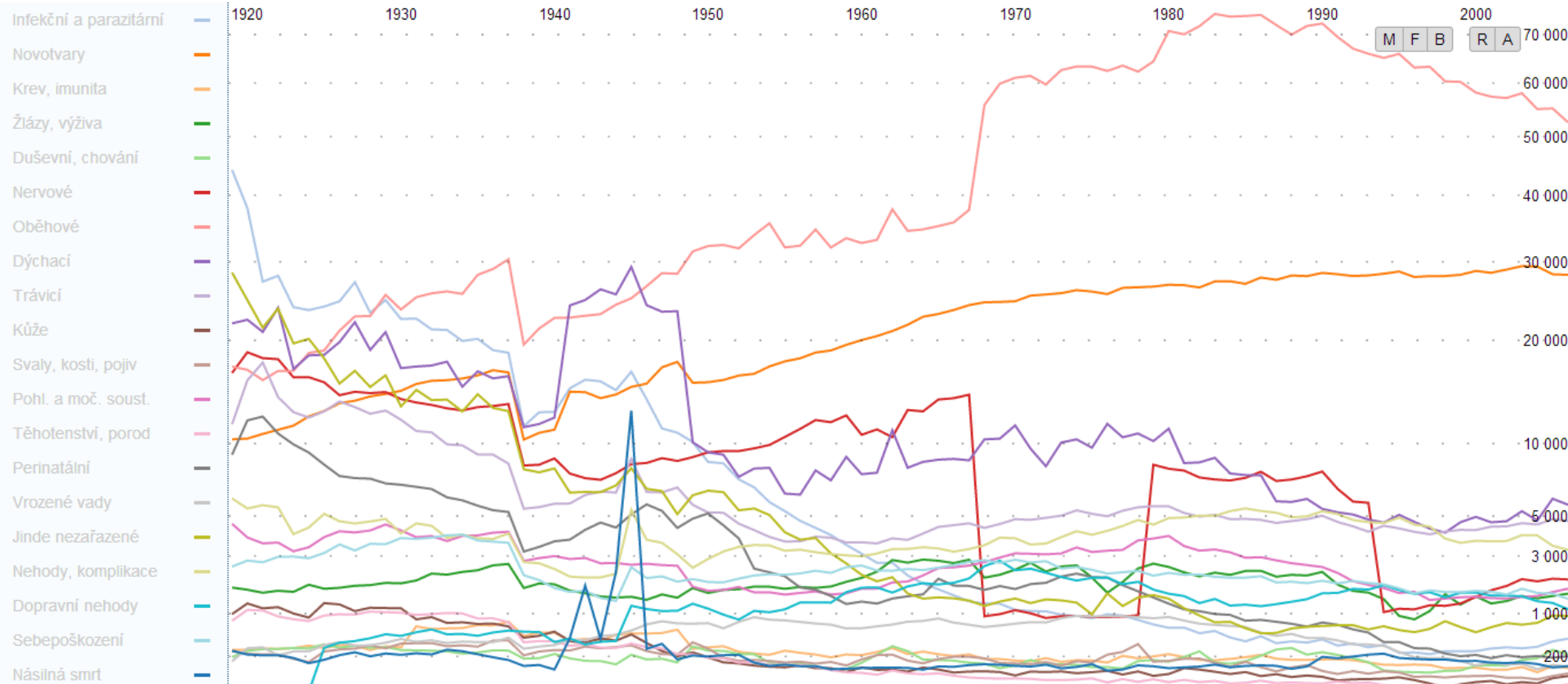
Připravená data
Standardní formáty
Málo práce
Weby institucí

Data neexistují
Existují, ale jsou
tajná
Spousta práce
(I programování!)

NUUDA × KRÁSA

Data dostupná všem

- ČSÚ (nuda, krása)
- Ministerstva, ŘSD, ÚZIS ...
- Globálně: Eurostat, OSN, WB, WHO ...
- Lokálně: Open Data (USA, jinde)
- Tipy: **Datablog**

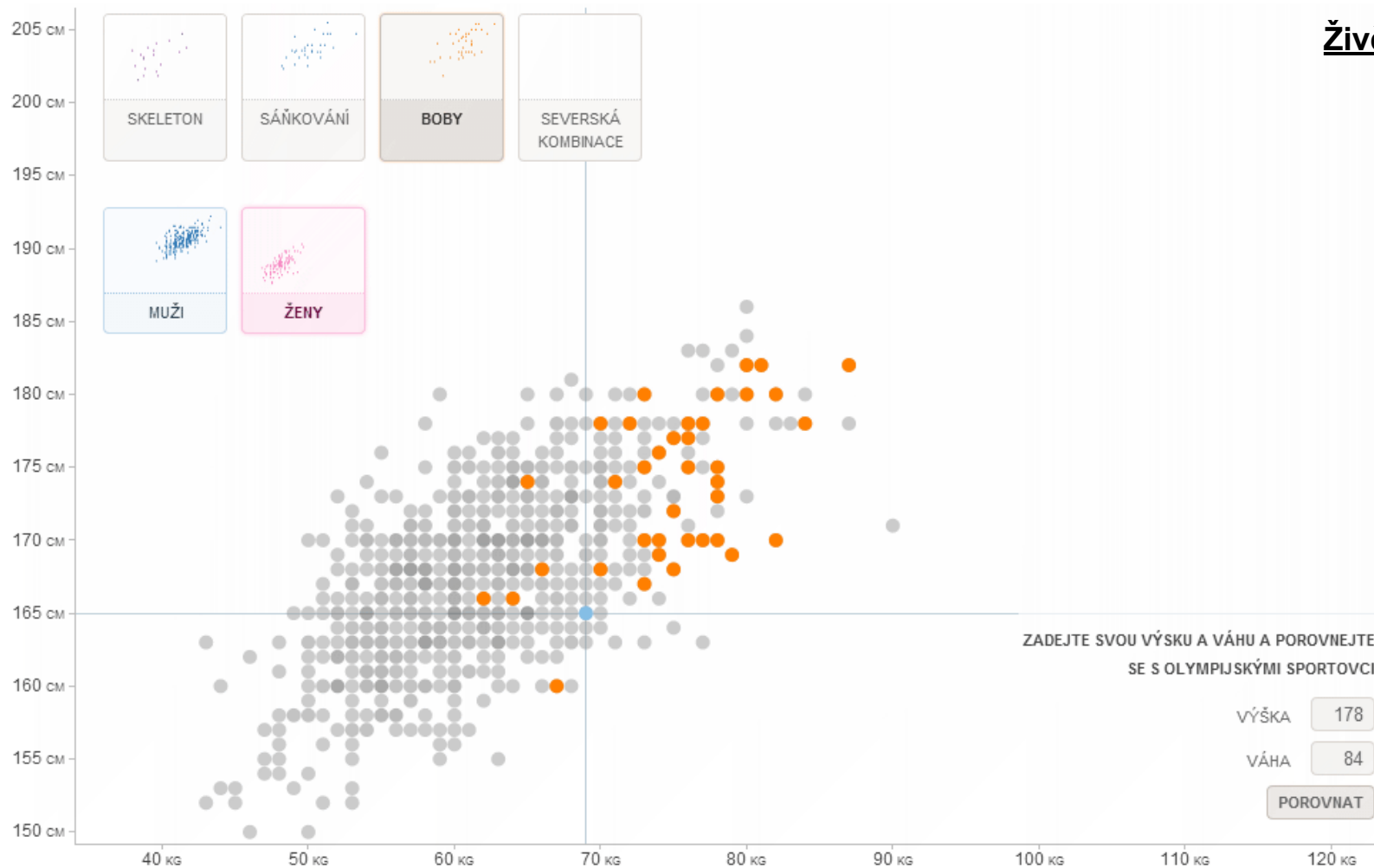


Data dostupná všem, ale

- Jak vypadají data pro volební prognózy ČT?
- Kolik platí VZP za konkrétní diagnózy?
- Kdo cestuje s prezidentskými delegacemi?
- **Odpověď: 106!**
- Pozor na zapadlé kouty webu, PDF

Data, co zatím neexistují

- Diskuze o mazání na Wikipedii (interaktivně)
- Barvy v různých kulturách (interaktivně)
- Vzorce brněnských semaforů
- **FF odpověď**: papír, nůžky, pastelky a ∞ času
- **FI odpověď**: scrapování



ZADEJTE SVOU VÝŠKU A VÁHU A POROVNEJTE SE S OLYMPIJSKÝMI SPORTOVCI

VÝŠKA 178

VÁHA 84

POROVNAT

Formáty souborů

- XLS



Formáty souborů

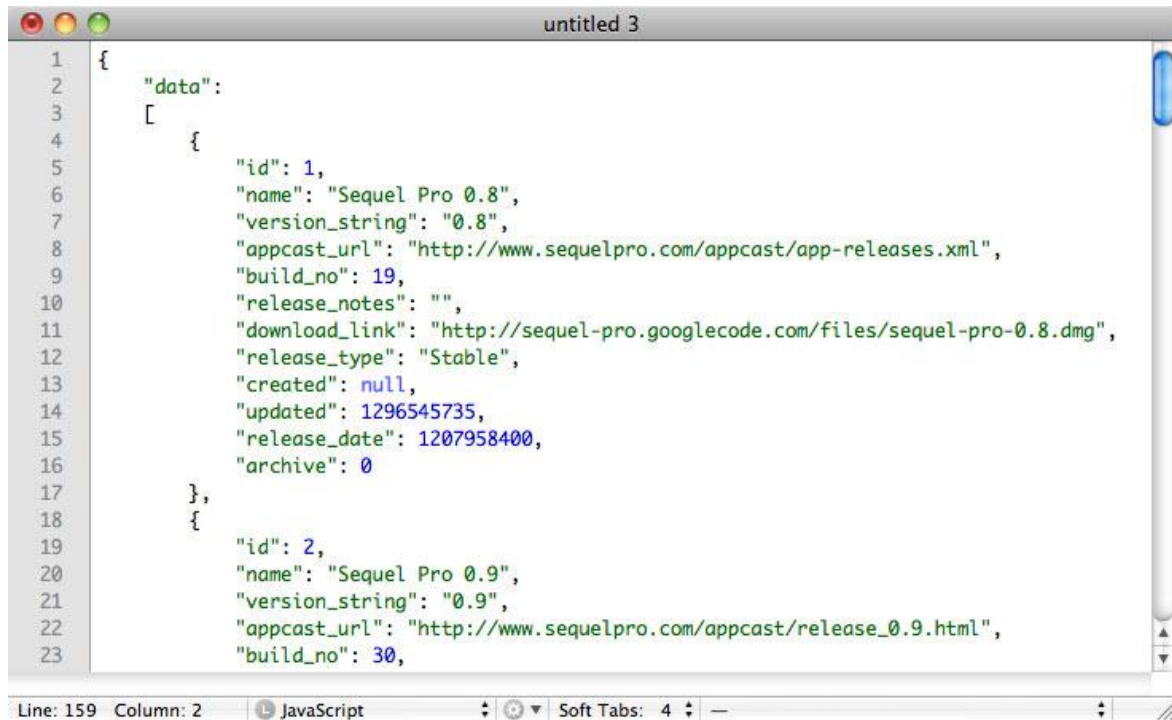
- XLS
- CSV



CSV is the data Kalashnikov: not pretty, but many wars have been fought with it and kids can use it.
#okfest

Formáty souborů

- XLS
- CSV
- JSON



```
untitled 3
1  {
2    "data":
3    [
4      {
5        "id": 1,
6        "name": "Sequel Pro 0.8",
7        "version_string": "0.8",
8        "appcast_url": "http://www.sequelpro.com/appcast/app-releases.xml",
9        "build_no": 19,
10       "release_notes": "",
11       "download_link": "http://sequel-pro.googlecode.com/files/sequel-pro-0.8.dmg",
12       "release_type": "Stable",
13       "created": null,
14       "updated": 1296545735,
15       "release_date": 1207958400,
16       "archive": 0
17     },
18     {
19       "id": 2,
20       "name": "Sequel Pro 0.9",
21       "version_string": "0.9",
22       "appcast_url": "http://www.sequelpro.com/appcast/release_0.9.html",
23       "build_no": 30,
```

Line: 159 Column: 2 JavaScript Soft Tabs: 4

Formáty souborů

- XLS
- CSV
- JSON
- **Mr. Data Converter**

Čištění dat



Čištění dat

- Vyčistit bordel (typicky: veřejné zakázky)
- Učesat klíčové proměnné
- Smazat ostatní
- Nejistota!

Google Docs

- Zkopírujte si data (sportovci v Soči/sochi.jdem.cz)
- Disciplína, stát, pohlaví: přeložit
- Jméno: rozdělit na křestní a příjmení
- Výška a váha: rozhodnout, kdo má zhubnout
- Výška a váha: orientační graf pro sporty