**Data Science in Python: Past, Present & Future**

**(PyData Track @ PyConLT 2019, SMK; #PyConLT19 #PyConLT)**

**25. 5. 2019**

**Radovan Kavický, GapData Institute & PyData Bratislava #PyDataBA**

# Talk Structure

- 1. **Introduction, History of Data Science & Why Python for Data Science?/Motivation** (5-8 min.)

- **Humble Beginnings of Data Science**

- **What is Data Science & Data Science Workflow**

- **History of Python & Future of Data Science**

- 2. **Main (tools), Where to start & How to do it?** (10-12 min.)

- **Python & Data Science (IDE-s and tools)**

- **Why should people care? (PyData Worldwide & PyData movement/educational programme)**

- 3. **PyData Bratislava (where we were/are & our activities), Partners, Meetups & Conferences** (2-5 min.)

- 4. **Conclusion (the end of talk +vision), Potential & Results** (2-5 min.)

- **Future of Data Science (in/outside Python)**

- **Data-Driven Approach. Everywhere. Now. (Citizen Data Science)**
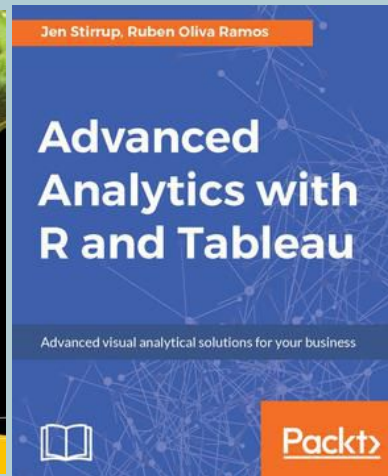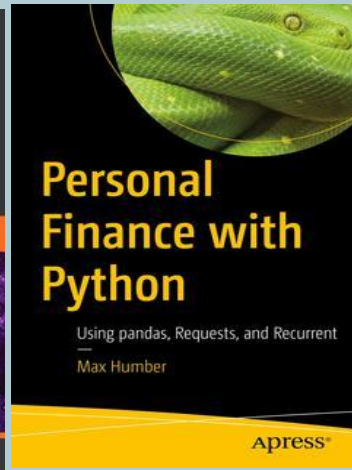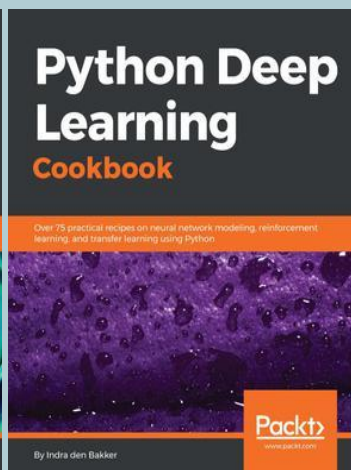
**Data Science in Python: Past, Present & Future**

# About me

- Economist (Macro, Finance, Statistics)
- Principal Data Scientist, Consulting (public, private)
- R, Python, Julia, Tableau > Matlab, SAS, Stata
- Data Science & Open Data & Public Policy
- PyData Bratislava, R <- Slovakia,  SlovakiaJUG, skczTUG football/soccer (FC Economist), Effective Altruism Slovakia

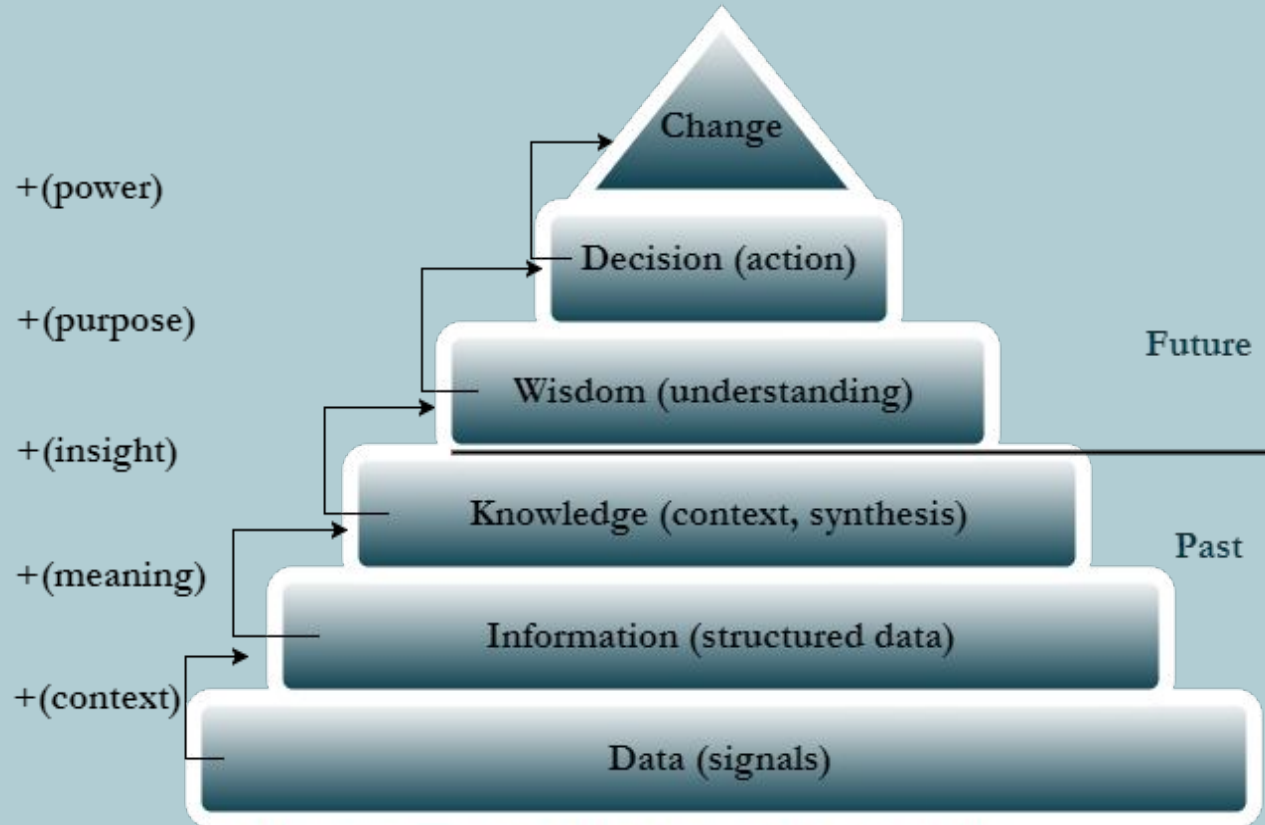Data Science in Python: Past, Present & Future

# How old is Data Science & who "coined" it?

- John Tukey, Bell Labs (+John von Neumann, 50's), DJ Patil (1st US Chief Data Scientist + Jeff Hammerbacher "coined term Data Science", 2012)

- collection of scientific results and methods for transformation of data from raw form to meaningful information, knowledge and wisdom, which should support better decisions

- Data Science ≠ Big Data

- Data Science: Statistics + programming; data analysis + Computer Science, modelling + Econometrics, Big Data, ML/DL

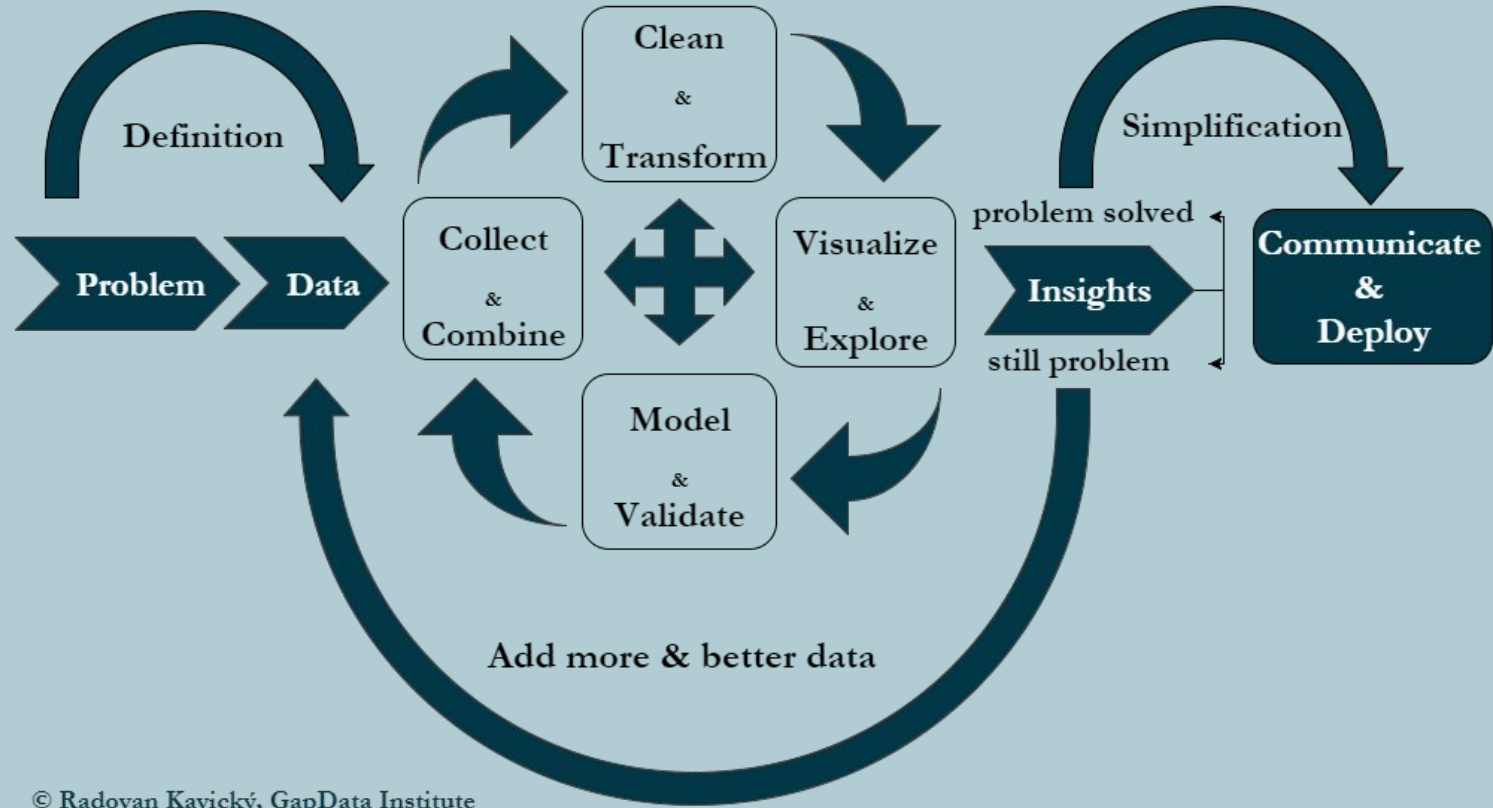- More on History of Data Science: Donoho, David (2015, Stanford), https://courses.csail.mit.edu/18.337/2015/docs/50Years DataScience.pdf

**Data Science in Python: Past, Present & Future**

# Data Science Pyramid



+(power)

+(purpose)

+(insight)

+(meaning)

+(context)

Change

Decision (action)

Wisdom (understanding)

Knowledge (context, synthesis)

Information (structured data)

Data (signals)

Future

Past

© Radovan Kavický, GapData Institute (inspired by DIKWD model by Curt Swindoll)

**Data Science in Python: Past, Present & Future**

Data Science Process & Workflow

© Radovan Kavický, GapData Institute

Data Science in Python: Past, Present & Future

# Where & When did Python



- 80's (idea), December 1989 (beginning), Guido van Rossum (Netherlands/Python is Dutch, CWI Research Institute)
- 1994 (1.0), 2000 (2.0), 2008 (3.0/3000/Py3K), 2020 (4.0 +end 2.0)
- hobby project during holidays
- Monty Python Flying Circus
- ABC language (cancelled +SETL), Amoeba (OS) needed script language
- First PC (ideal), universal a multiplatform+ easy to use, hard to master
- Not for data analysis and Data Science (s-l-o-w)
- IPython/Jupyter Notebook (Fernando Perez) –JUlia, PYthon, teR

**Data Science in Python: Past, Present & Future**

# R & Data Science



- **Programming language S - John Chambers, Bell Labs (USA)**

**+ Rick Becker, Allan Wilks**

**1976, Fortran subroutines, later S-Plus (commercial version)**

- **R - Robert Gentleman a Ross Ihaka (New Zealand, University of Auckland)**

**August 1993, C + Fortran**

- **R as implementation of S (didn't want to pay license)**

**Data Science in Python: Past, Present & Future**

# Python & Data Science (tools)

**IDE's – Jupyter Notebooks (IPython kernel), PyCharm, JupyterLab**

**Data Collection – Feather (binary file format/non-csv/Apache Arrow/Wes McKinney), pandas datareader**

**Data Visualization – Seaborn (matplotlib based/static), Bokeh (interactive/d3.js like), Plotly (declarative dataviz), Altair (static/js)**
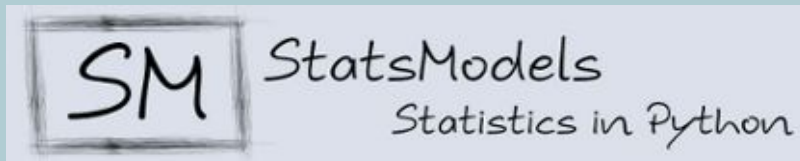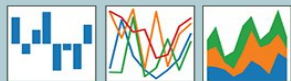
**Data Cleaning & Transform - datacleaner (automate cleaning your data in Pandas), Blaze (NumPy/pandas-like), Dask (parallel computing)**

**Data Modeling – StatsModels +Patsy (describe statistical models), PyMC3 (Bayes/statistical modeling), Keras (TensorFlow & PyTorch/DL)**

pandas
$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$

SM StatsModels
Statistics in Python

PyMC3

**Data Science in Python: Past, Present & Future**

# R vs/& Python "war" in Data Science



R - Studio (IDE for R & Python)

Jupyter Notebook (Irkernel) + r-essentials (conda)

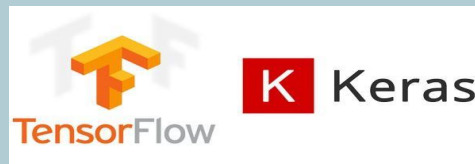Anaconda (open-source distribution of Python&R)

Tidyverse (tidyverse.tidyverse.org, Magrittr, %>%)

dplyr, tidyr (subsetting, joining, data manipulation)

ggplot2, rBokeh, ggvis (data visualization)

caret – modelling (regression, classification)

H2O.ai, Tensorflow & Keras – ML & DL

Shiny, blogdown, R Markdown (publication, RMd, blog)



**Data Science in Python: Past, Present & Future**

# Future of Data Science (Pythonic or Future is Female/Julia?)





- **Fast as C, syntax/easy to read as Python, dynamic as Ruby, mathematic notation as Matlab, statistics as in R, string processing as Perl, true macros like Lisp**
- **Creators (2009, 2012, MIT): Viral B. Shah, Stefan Karpinski, Jeff Bezanson, prof. Alan Edelman**
- **Julia Computing (developers 1.0 – summer 2018)**
- **Future of Data Science & designed with Data Science in mind**
- **Python not for analysis and Data Science (s-l-o-w)**
- **Julia is (f-a-s-t/ LLVM compiler)**
- **Solves "two-language problem" - prototype (slow, dynamic), delivery (quick, static)**
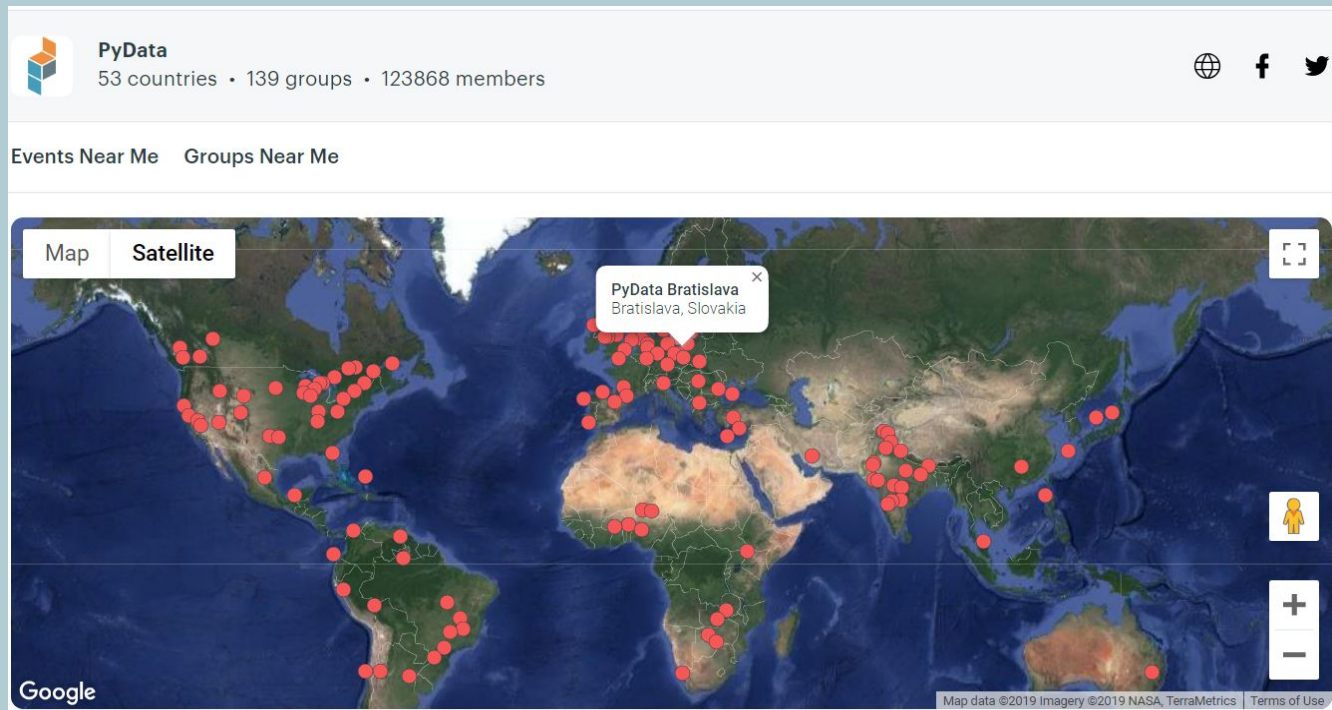- **IPython/Jupyter (Fernando Perez) –JUlia, PYthon, teR**

**Data Science in Python: Past, Present & Future**

# PyData as global movement



PyData
53 countries · 139 groups · 123868 members

Events Near Me    Groups Near Me

PyData Bratislava
Bratislava, Slovakia

Source: https://www.meetup.com/pro/pydata/

- **More Meetups outside US than within**
- **Slovakia (40th globally):**
- **Few days after us LA (as 41st & 11th in US)**
- **Before us: Poland (Warsaw, Wroclaw, Krakov), Austria (Vienna, few months)**
- **After us: Czechia (Prague, 2 years), Hungary (Budapest, 1 year)**

**Data Science in Python: Past, Present & Future**

NUMF⊙CUS
OPEN CODE = BETTER SCIENCE

April 2019 Newsletter

**The Headlines:**

The recent breakthrough image of a black hole was powered by multiple NumFOCUS projects!

## Dr. Katherine Bouman

- Numpy (van der Walt et al. 2011)
- Scipy (Jones et al. 2001)
- Pandas (McKinney 2010)
- Astropy (The Astropy Collaboration et al. 2013, 2018)
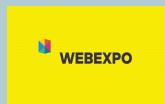- Jupyter (Kluyver et al. 2016)
- Matplotlib (Hunter 2007)

**Data Science in Python: Past, Present & Future**

# PyData Bratislava & activities within CEE + V4



- 300+ members
- 600+ on mailing list
- 800+ followers
- Partners: O'Reilly, robime.it, Learn2Code, sli.do
- Sponsors: H2O.ai, SwissRe, ERNI, Kiwi Security
- WeAreDevelopers World Congress ("GAPDATA-25" & "PYDATA-25"), Berlin, 6.-7. June
- OSCON 2019 ("PCGAPDATA" for 25% discount)
- TechSummit, Bratislava, 29.-30. May ("pydatabratislava20" for 20% discount)
- Our Conferences (satRday, PyData Bratislava)

**Data Science in Python: Past, Present & Future**

# #AllForJan Hackathon (cooperation with FIIT STU)



Source: https://github.com/AllForJan

- Data journalism
- Public Policy, Activism & Data
- Useful tools for journalist
- Agricultural data (dotations)
- Vice-dean FIIT Pištek & GapData Institute
- LOD (linked open data), transparency check
- All outputs public (GitHub)
- Over 50 developers & data scientists

Data Science in Python: Past, Present & Future

# GapData Institute (GDI) and how to support us.

⧠ **Economic Research & Public Policy & Data Science think-tank (data-tank)**

⧠ **Data. Think. Change.**

⧠ **GapData Institute (GDI) is a non-profit nonpartisan research institution harnessing power of data & wisdom of economics for public good.**

⧠ **Transparent account (from day #1; SK7383300000002200933920 https://www.fio.sk/ib2/transparent?a=2200933920)**

⧠ **Partnership (openness, transparency)**

⧠ **Slides (this talk): tiny.cc/pyconLT2019vilnius**

⧠ **https://github.com/radovankavicky/PyConLT2019**

**GAP DATA INSTITUTE**

**NUMFOCUS**
**OPEN CODE = BETTER SCIENCE**

**PyData**
*Bratislava*

PAY by square

**Data Science in Python: Past, Present & Future**

#PyDataBLN    **Data Science in Python: Past, Present & Future**

# Thank you for your attention

**Contact:**

**Radovan Kavicky**

radovan.kavicky@gapdata.org

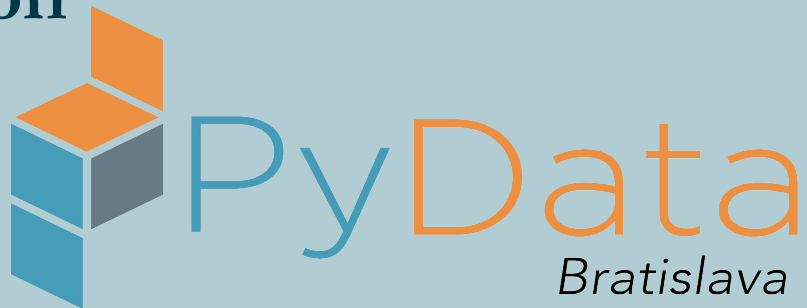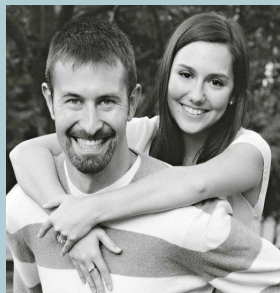radovan.kavicky@gmail.com

**+421 949 716 214 (SK)**

http://www.linkedin.com/in/radovankavicky

https://gapdata.slack.com/

Invitations (e-mail): https://gapdata.herokuapp.com/
https://github.com/radovankavicky
https://github.com/GapData/PyDataBratislava

@radovankavicky, @PyDataBA, @GapDataInst

**PyData** *Bratislava*

**GAP DATA INSTITUTE**

In case you have any question, feel free to ask.