# Artificial Intelligence

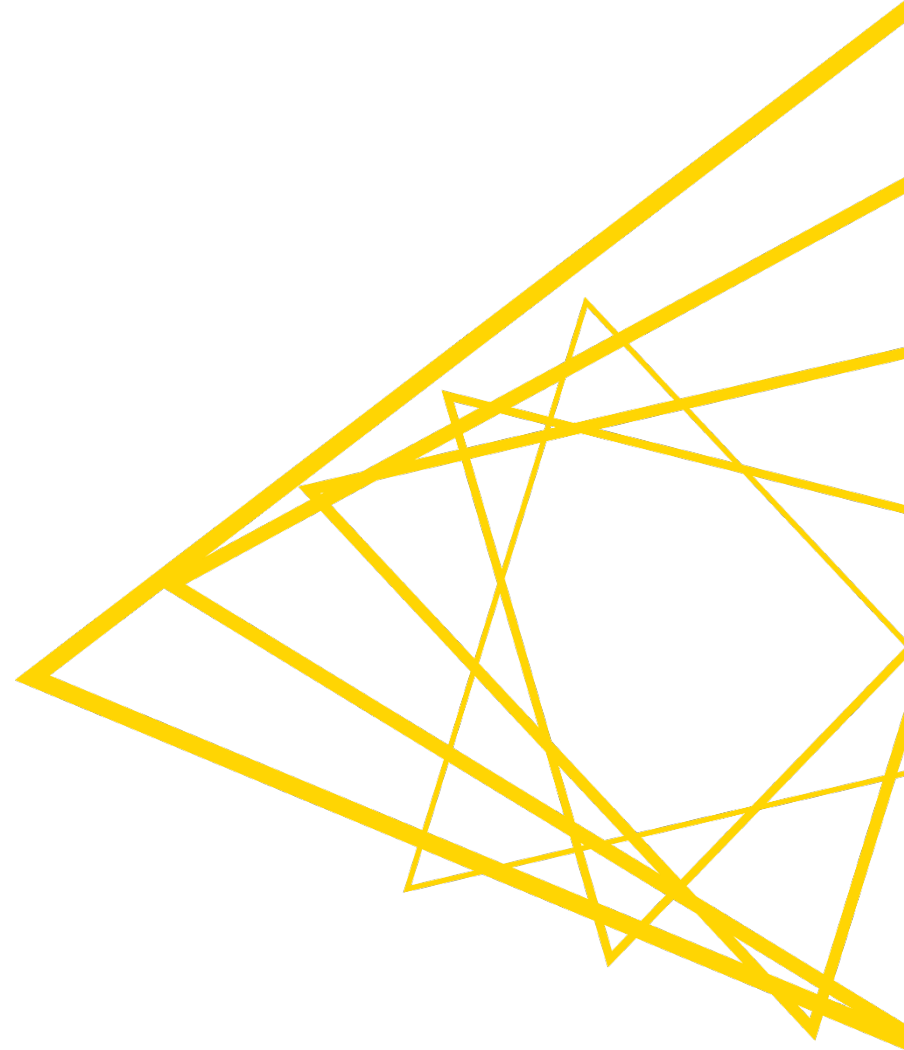# Visual Programming for ~~Data Science~~

Working with Data at the right Level of Abstraction

**Prof. Dr. Michael Berthold**
ISACA Fokus Event, Nov 10, 2022

# Data Science Personas

# Data Science Researchers

Researchers focus on Data Science algorithms

- often have CS/NumMath/Stats background

- use Python, Javascript, R, SQL, C, Java, or XYZ

- tend to not worry about others.

# Data Scientists

Data Scientists

- have (also) a specific domain background
- often can program in Python, R, SQL, C, Java, or ...
- …but really focus on the data flow
- want to use and try out other stuff

# Casual Users ("Citizen Data Scientists")

The occasional data person

- wants to start from something existing that solves a similar problem

- often comes from an Excel, BI, or other data world

- needs solid documentation and proper abstraction

# Data Science Consumers

Data Science Consumers use

- their favorite frontend
(= Data Science needs to embedded via an API)

or

- Interactive Analytics (= Data Science Applications)

# Data Science Personas



Coders:
Programs & Scripts

Data Scientists:
Visual Programming

Citizen Data Scientists:
Components & Blueprints

Data Science Consumers:
Applications & Services

Data Science
Algorithm Inventors

Data Science
Creators

Data Science
Consumers

# No Code, Low Code, Visual Programming?

# No Code

No Code for Simple Problems:

| No Code | (Code-Based) Programming |
|---------|--------------------------|

→ Complexity



Visual Environment allows quick creation of solutions for standard problems.

# Low Code

Low Code for Standard Problems:



| No Code | Low Code |
| | (Code-Based) Programming |

Complexity

Visual Environment allows 80% creation of solutions.

…for the rest (& details): reach out to code.

(often: Visual UI on top of **one** programming language.)

KNIME
Open for Innovation

# Visual Programming

Visual Programming for all types of Problems:



Visual Programming

(Code-Based) Programming

Complexity

Visual Environment allows complete creation of solutions.
(from simple to complex)

If desired: you can also include code (of various languages!).

(Visual UI models what's running under the hood.)

# Visual Programming
# for Data Science

# The Data Science Life Cycle

**Blend & Transform**

**Validate & Deploy**

**Model & Visualize**

**Creation**

**Production Process**

**Production**

**Consume & Interact**

**Optimize & Capture**

**Monitor & Update**

Open for Innovation
KNIME

# Data Science Activities



Diagnostic Analytics
(Creating Insights)

Descriptive Analytics
(Automating Repetition)

Advanced Analytics (AI/ML…)
(Predictive, Prescriptive)

Data Wrangling
(Virtual Warehouses)

Productionize
(Deploy, Manage, Govern)

KNIME
Open for Innovation

# No Code for Analytics

# No Code for Analytics

Benefits:

- (Reliable) Automation & Reproducibility
- Documentation

# Creating Data Insights: Abstraction helps

ARR Cohort Analysis for Excel Users

# Creating Data Insights: Abstraction helps

ARR Cohort Analysis for Excel Users – No Code Encapsulated:

# Abstraction in KNIME: Components

# Low Code for Data Wrangling

# Flexible Data Wrangling

# Flexible Data Wrangling

# Low Code for Data Science

# No Code for Data Science

# No Code for Interactive Visualizations

# Codeless Deep Learning

# Data Science *is* complex: Visual Programming

Visual Data Science Algorithms: Recursive Feature Elimination

"Computer programming is an art, because it applies accumulated knowledge to the world, because it requires skill and ingenuity, and especially because it produces objects of beauty."

*Donald Knuth 1974*

Open for Innovation
KNIME

# Low Code for Data Science

Leveraging Code in Visual Programming:



…works with (and from…) Jupyter Notebooks as well.

# Data Science *is* complex: More Topics

Data Science "Stuff":

- Data Blending:
  - Feature Selection
  - Feature Engineering
  - Anonymization
  - …

- Modeling
  - Model Selection
  - Parameter Optimization
  - Ensemble Creation
  - Active Learning
  - Transfer Learning
  - …

Management Requirements:

- Compliance
  - Best Practices
  - Required Standards
  - Bias Detection and Removal
  - Data Privacy / GDPR
  - Explainable "AI"
  - …

- Governance
  - Integration / Dependencies
  - Traceable / Lineage
  - Documented
  - Reproducibility
  - …

KNIME
Open for Innovation

# Data Science Teams

# The Data Science Life Cycle

**Blend & Transform**

**Model & Visualize**

**Deploy & Manage**

**Consume & Interact**

**Techniques:**

- Joins
- Pivot
- Normalization
- PCA
- …

- Clustering
- Regression
- Deep Learning
- Visualizations
- …

- Model Monitoring
- Updating
- Validation
- …

- Active Learning
- Interactive Applications
- Services
- …

**Technologies and Languages:**

- SQL
- Data Lakes (AWS, Azure, Google…)
- …

- R, Python
- Java, C
- JavaScript, …

- Go
- Kubernetes
- Docker
- Grafana
- …

- Plotly, R Shiny
- JavaScript
- Vue.js
- REST
- …

Data Science Teams shouldn't really have to worry about all of this…

Open for Innovation
KNIME

# Democratizing Data Science

What's the Essence of the Job:

- Create (complex) Data Science processes, often together with other Experts

What's usually not part of the Job:

- Inventing, Writing, Optimizing new Algorithms

What shouldn't be part of the Job:

- Caring about the details of the underlying technology
- Worrying about interfaces of different tools
- Worrying about library versions and (backwards) compatibility
- Manually doing repetitive tasks

KNIME
Open for Innovation

# Visual Programming for Data Science

# Technologies & Languages under the Hood

# Low Code for Deployment

# No Code Deployment of Data Apps

# Deploying Data Science as Service

# Deploying Data Science as Service

# Visual Programming for all Aspects of Data Science

Data Wrangling
(Virtual Warehouses)

Descriptive Analytics
(Automating Repetition)

Diagnostic Analytics
(Creating Insights)

Advanced Analytics (AI/ML…)
(Predictive, Prescriptive)

Productionize
(Deploy, Manage, Govern)

# Visual Programming for Data Science makes Sense!

- Data Science is a Team Sport:
  - requires a broad spectra of skills – often done in collaboration
  - requires understanding and access to the inner wheels of an algorithm (but not to the algorithm itself)
  - No single language or tool has all the answers – it's about tool blending as well!

- No-Code, Low-Code, Visual Programming (if done right):
  - automates repetitive tasks
  - provides the appropriate level of abstraction to allow focus on modeling a data process
  - allows access to all relevant nuts and bolts
  - allows to abstract / encapsulate sophistication for others to safely reuse
  - removes tool interaction/interface complexity
  - provides transparency & governance
  - (and it's not just a UI slapped on top of a language)

  - …allows everybody to gradually increase their Data Science Proficiency (in the same Environment).

Open for Innovation
KNIME

# Data Science Upskilling



Individual Sophistication

ML & AI

Stats and Models

Data Exploration

Automated ETL

Spreadsheets & BI

# Organizational Upskilling



Blend & Transform

Validate & Deploy

Model & Visualize

Creation

Production Package

Production

Consume & Interact

Optimize & Capture

Monitor & Update

Read

Transform

Analyze

**Organizational Maturity**

Explorative DS

Manual Deployment

Monitoring

Continuous Deployment

Open for Innovation
KNIME

# Upskilling Individuals and Organizations

# Want to get started?

# Getting Started?



Beginners corner

**[L1-DS] KNIME Analytics Platform for Data Scientists: Basics**
- Content 38 Modules
- Difficulty Basic
- Course Length 8 hours
- Rating ⭐⭐⭐⭐⭐ 1079 Reviews
- Reviewer/Instructor **Maarit Widmann**

Start

Welcome to the **KNIME Hub**

**First Steps to Building Workflows**

Example-based learning as a beginner can help you acquire new skills faster. Experts from the entire community share their workflows here for you to download and inspect. We collected a dedicated set developed by KNIME in the Beginners space below. These workflows are accompanied by a beginner 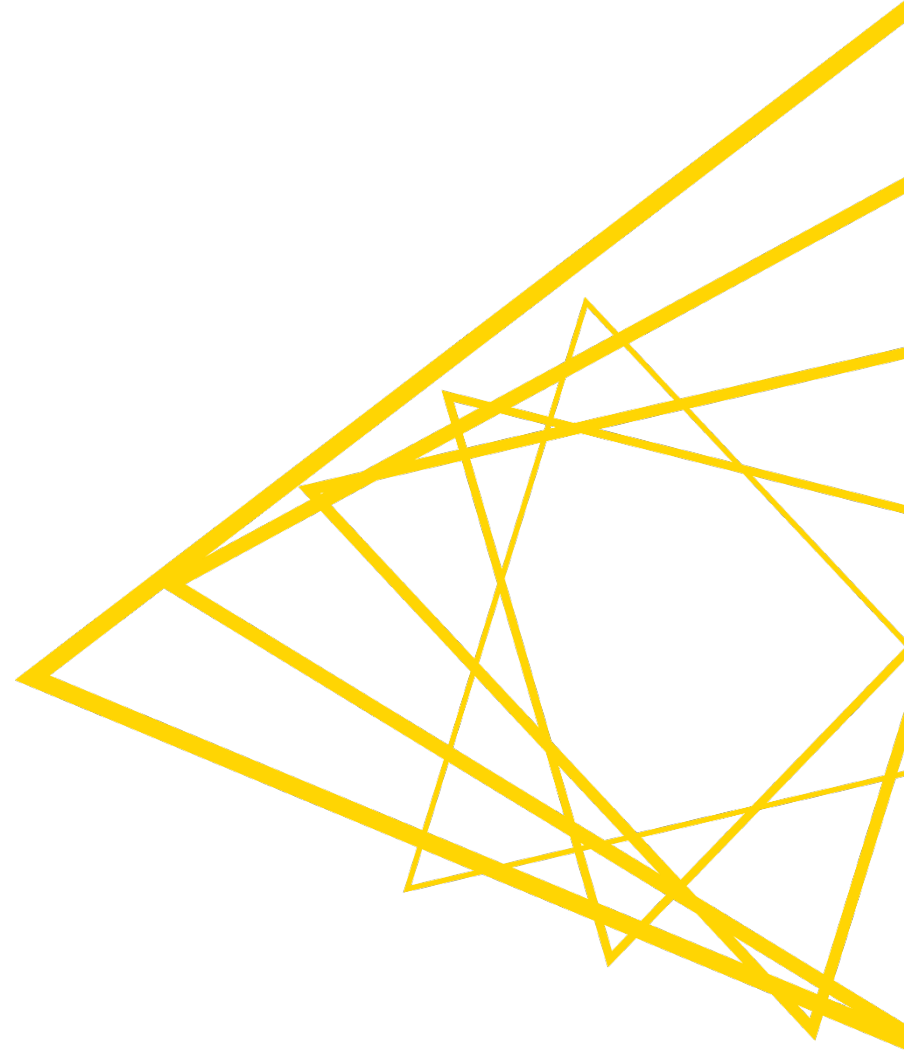cheat sheet and take you through your first steps to read, explore, transform, and deploy data in KNIME. Explore the links below.

**Explore relevant resources**
- The Beginners space on the KNIME Hub
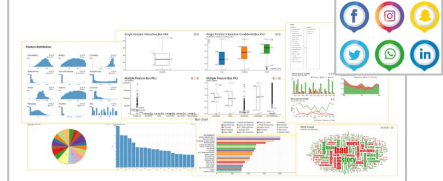- A cheat sheet on building workflows for beginners
- Getting Started Guide

**Data Exploration in #66daysofdata with KNIME**

A roadmap to deepen your knowledge of data exploration techniques within the #66daysofdata challenge

Co-author: _Rosaria Silipo_

Have you heard about the _#66daysofdata challenge_?

Certified
**Basic Proficiency**
in KNIME Analytics Platform

## https://www.knime.com/learning

## KNIME for Beginners

**Just KNIME It!**
Prove your KNIME knowledge and practice your workflow building skills by solving our weekly challenges.

**Basic KNIME**
31 videos • 108,603 views • Updated today

Here you can find videos about basic issues KNIME, like hiliting, file reading, and so on.

**Low Code for Advanced**
Home  About  Editor's Picks  Getting Started  Theory

Getting Started in Low Code for Advanced Data Science

Cheat Sheet: Building a KNIME Workflow for Beginners

**KNIME® BEGINNER'S LUCK**

A Guide to KNIME Analytics Platform for Beginners

Open for Innovation
**KNIME**