A complex network graph composed of numerous blue and white circular nodes connected by thin blue lines, forming a dense web-like structure.

Data Engineering + MLops SS24

Week 1 - Introduction

Prof. Dr.-Ing. Janis Keuper



INSTITUTE FOR MACHINE
LEARNING AND ANALYTICS



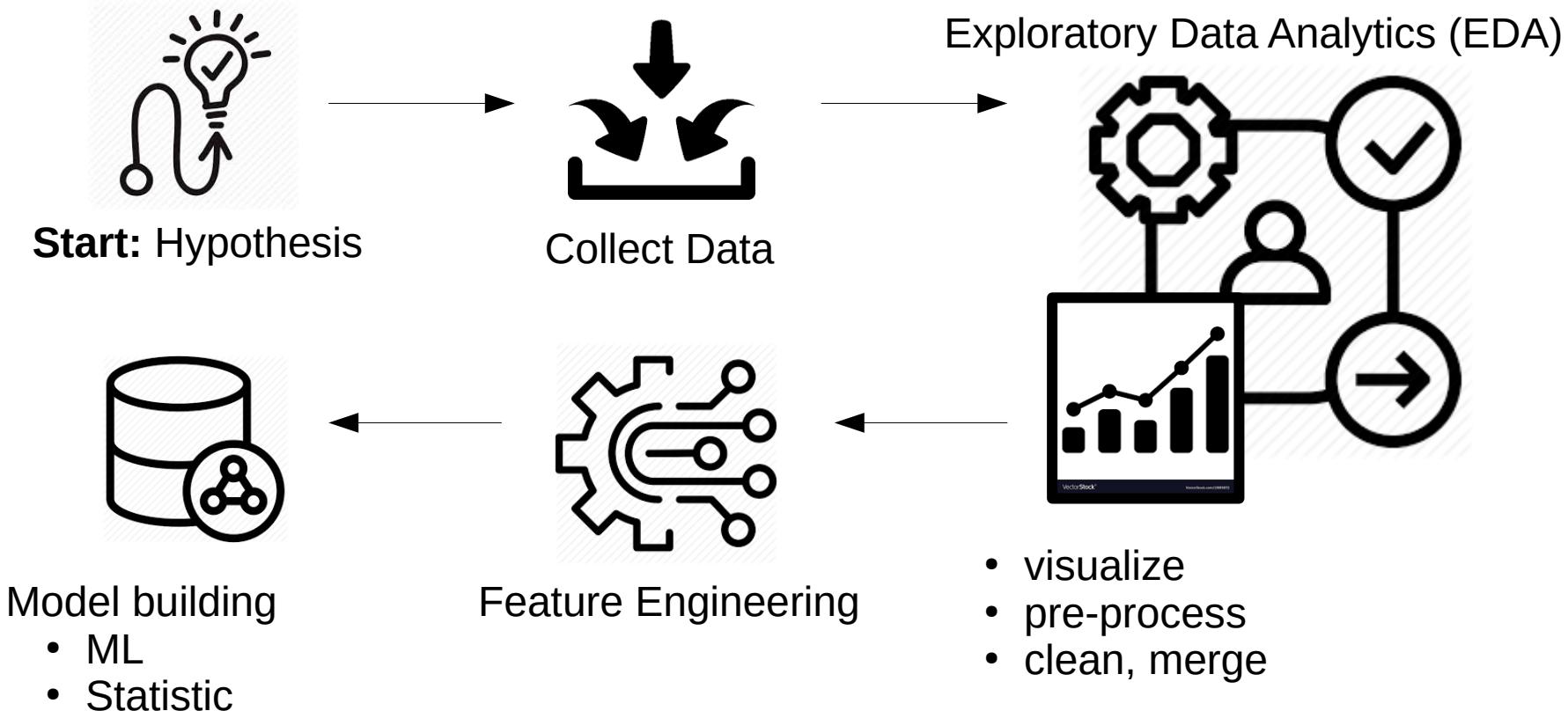
Hochschule Offenburg
offenburg.university



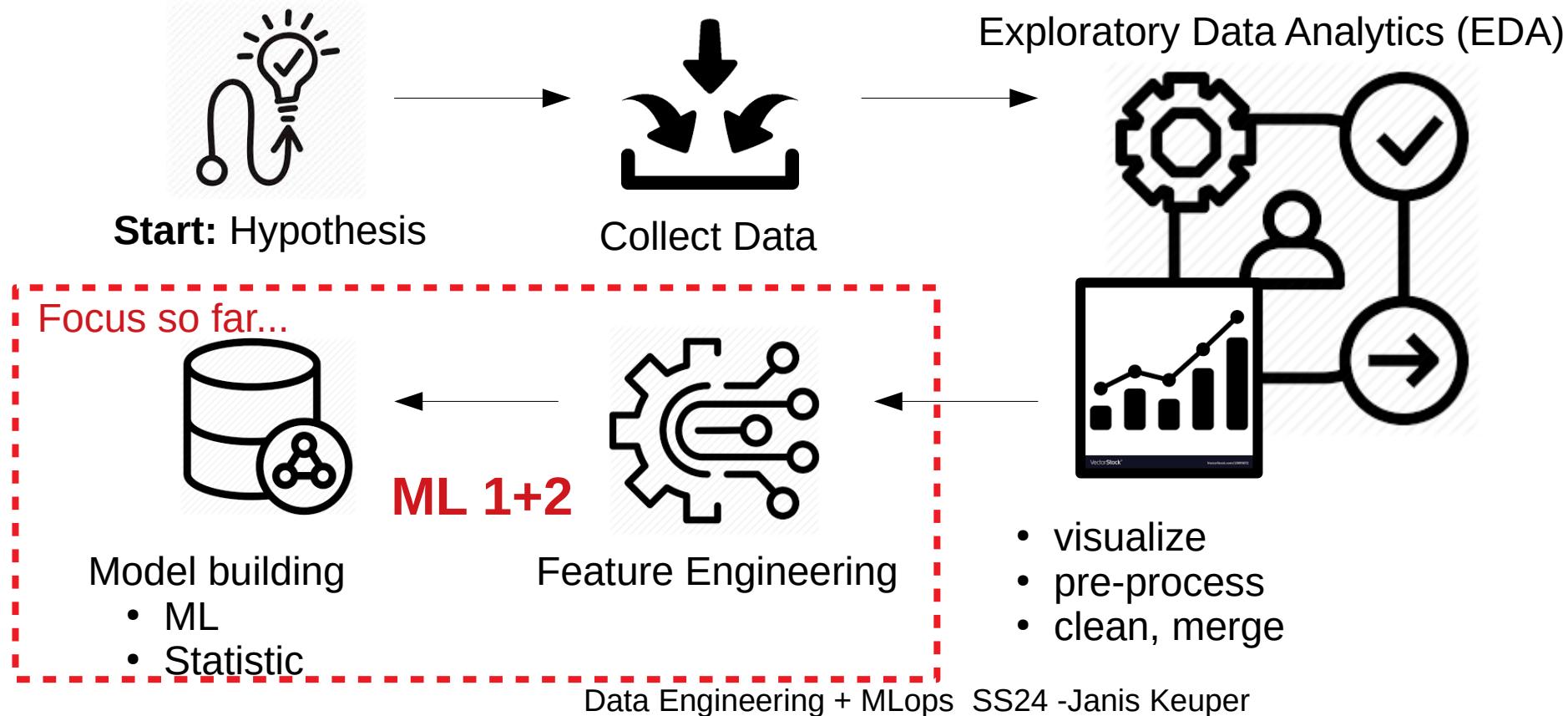
- Machine Learning Work Flow
- What is “Data Engineering ?”
- What is „MLOps ?“
- Python Data Ecosystem



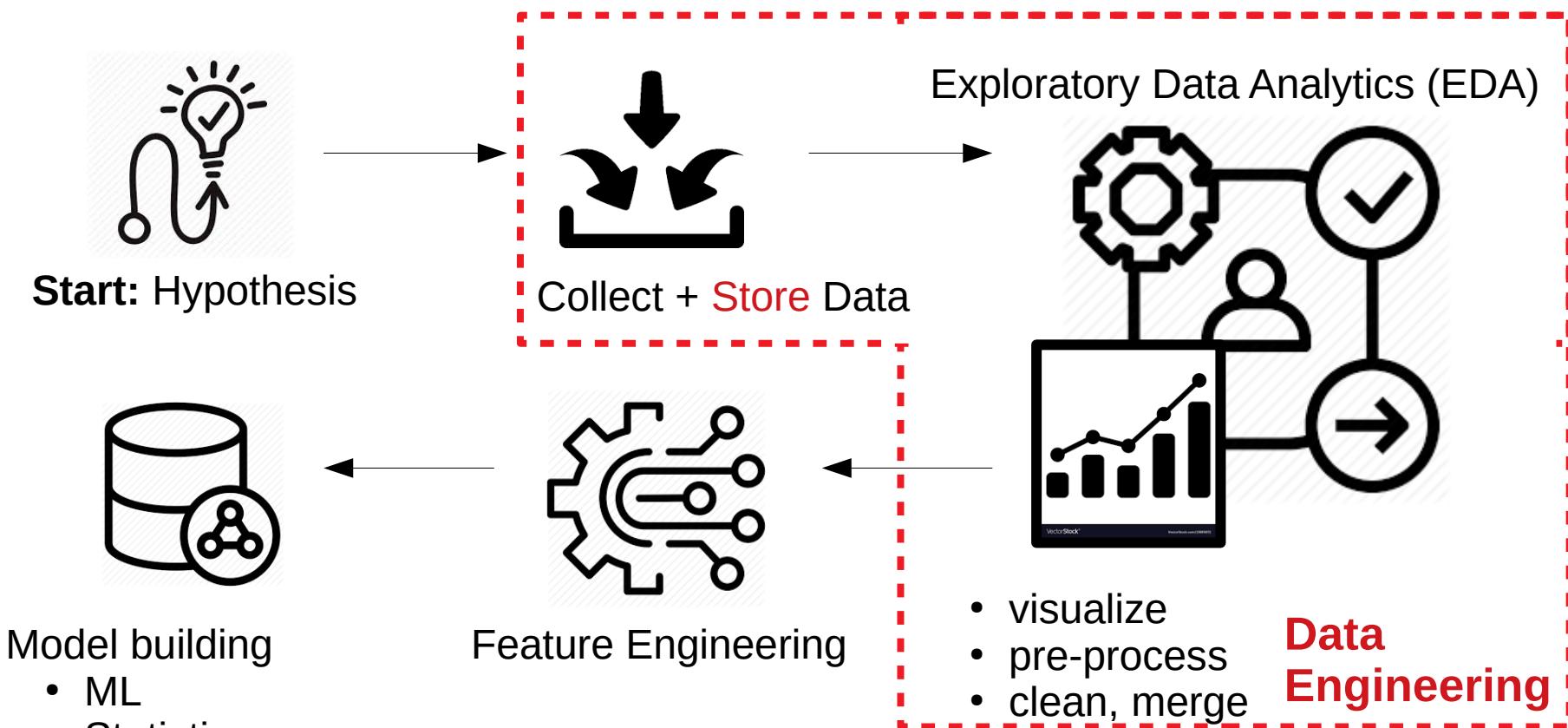
Simple Question: What is Data Engineering?

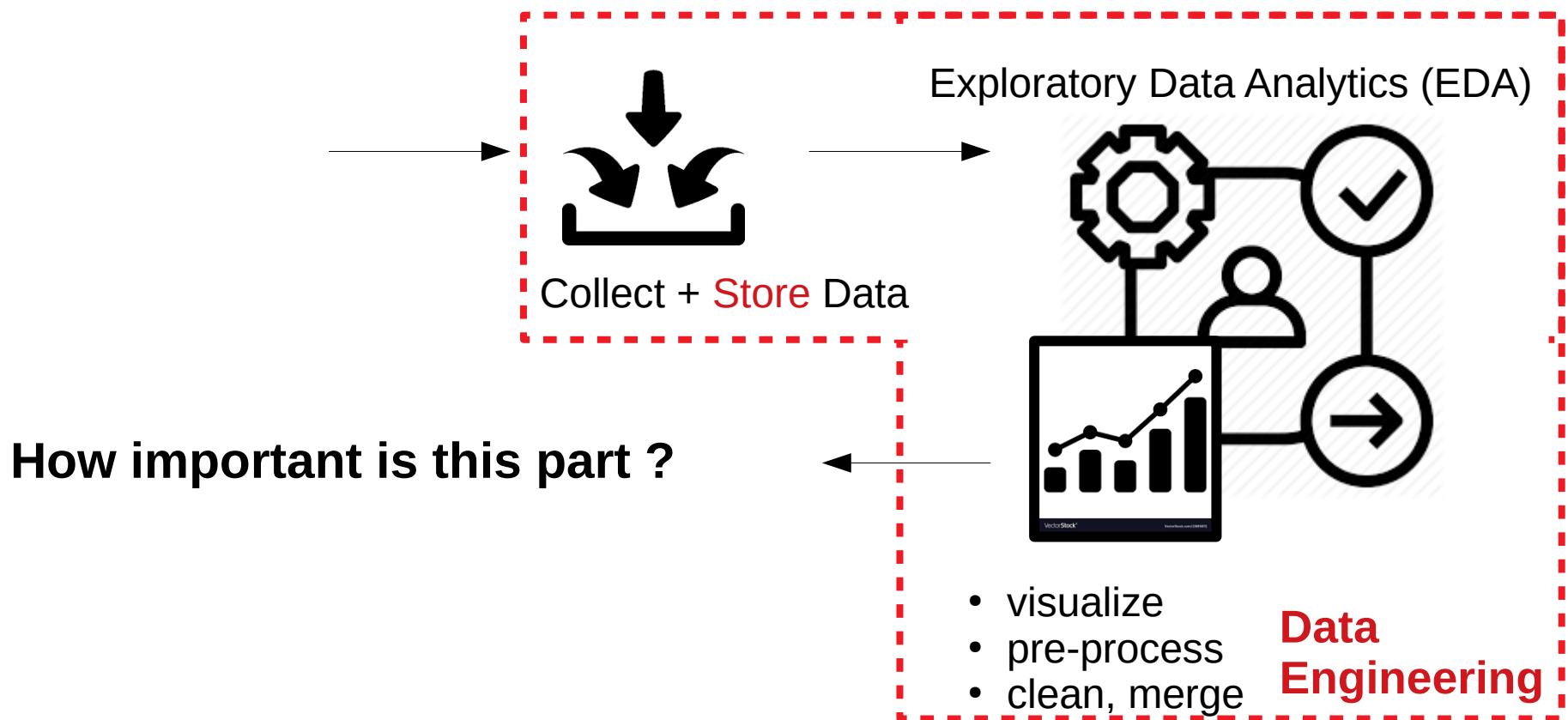


Data Science Workflow



Data Science Workflow





Simple Question: What is Data Engineering?



And how does it relate to or differ from... ?

Data Science

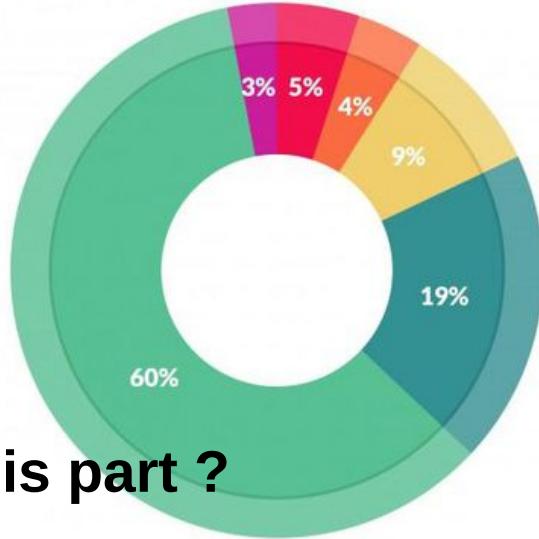
Big Data

Data Bases

Data Mining

Artificial
Intelligence

Machine Learning?



What data scientists spend the most time doing

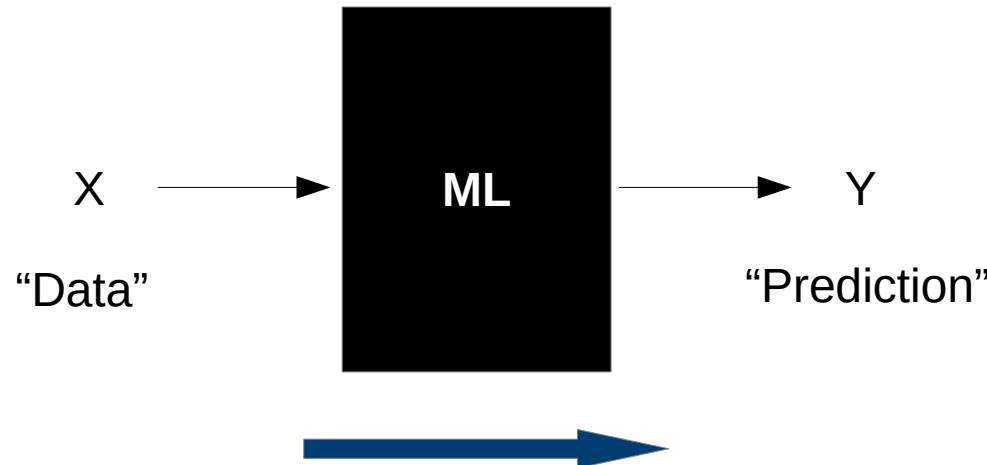
- Building training sets: 3%
 - Cleaning and organizing data: 60%
 - Collecting data sets; 19%
 - Mining data for patterns: 9%
 - Refining algorithms: 4%
 - Other: 5%
- ML

How important is this part ?

source: study by forbes.com [2]

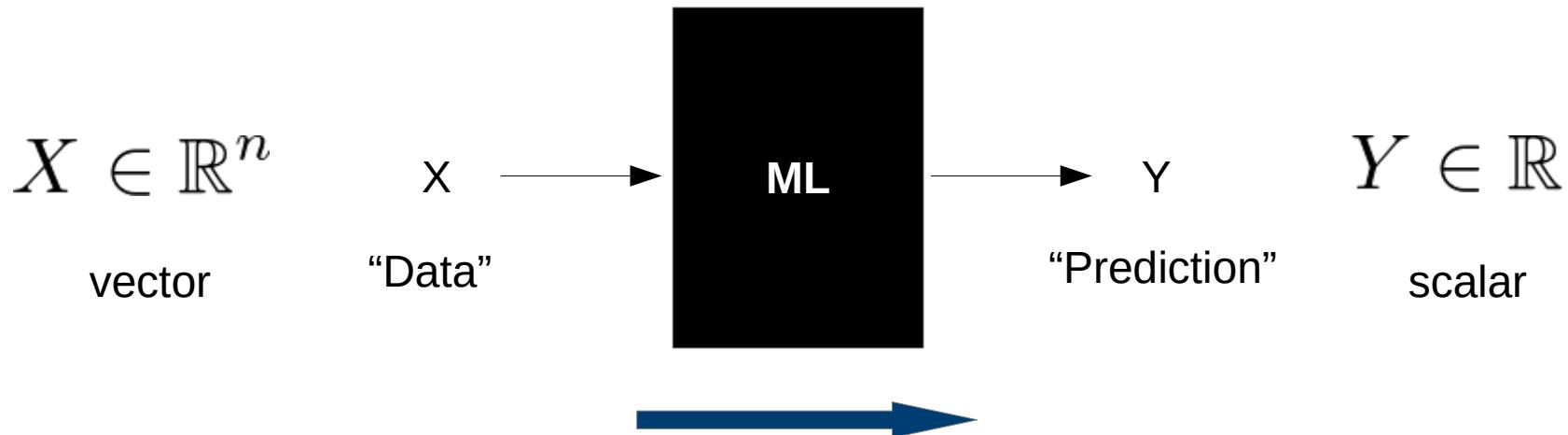
How important is this part ?

General abstract scheme of ML algorithms:



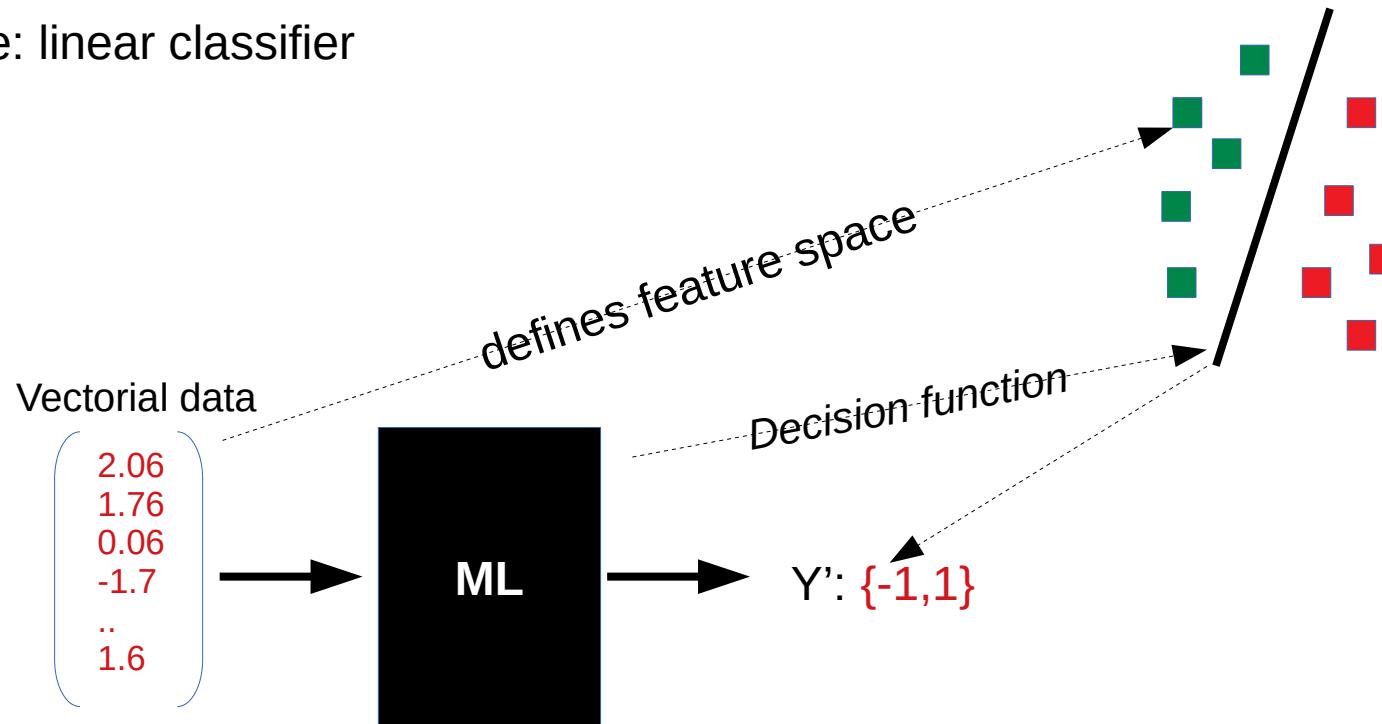
ML algorithms “learns” ***mapping*** from input to output by example data

General abstract scheme of ML algorithms (**so far**):

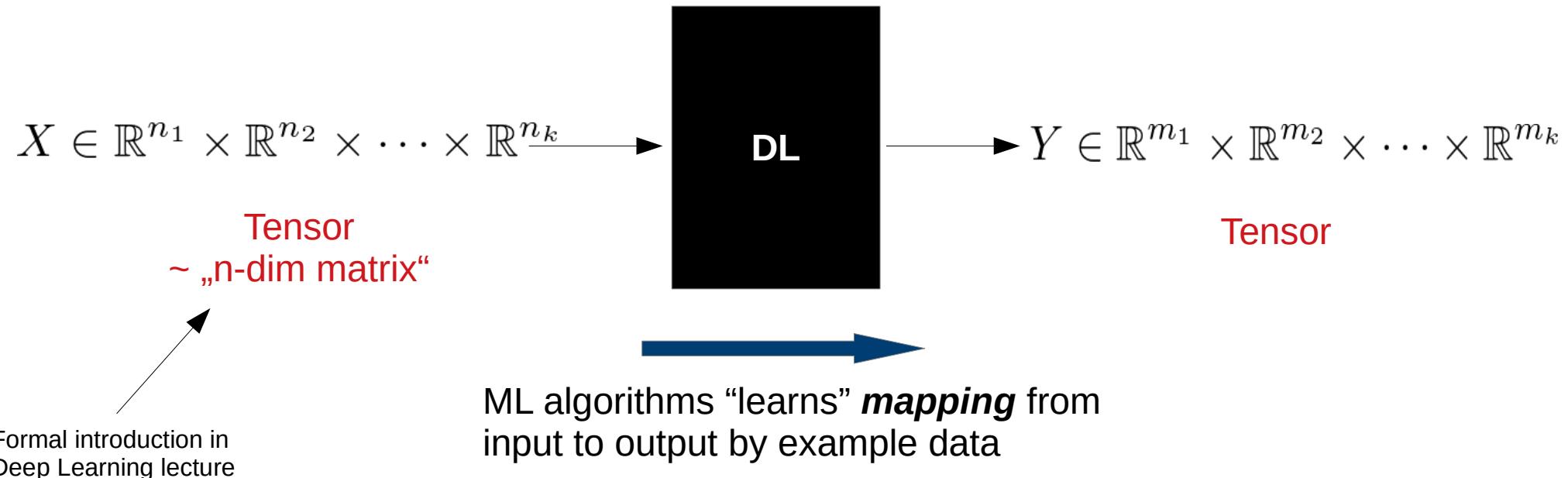


ML algorithms “learns” ***mapping*** from
input to output by example data

Example: linear classifier

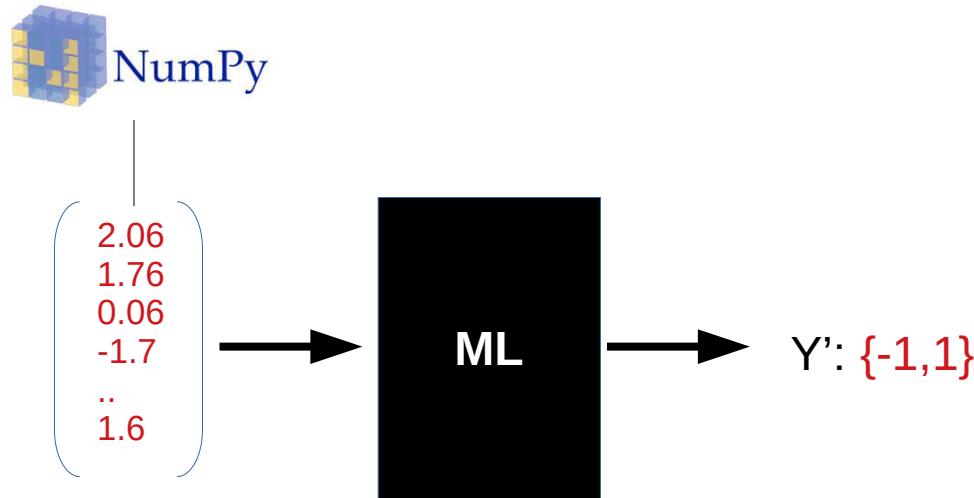


General abstract scheme of ML algorithms (modern algorithms, i.e. Deep Learning):



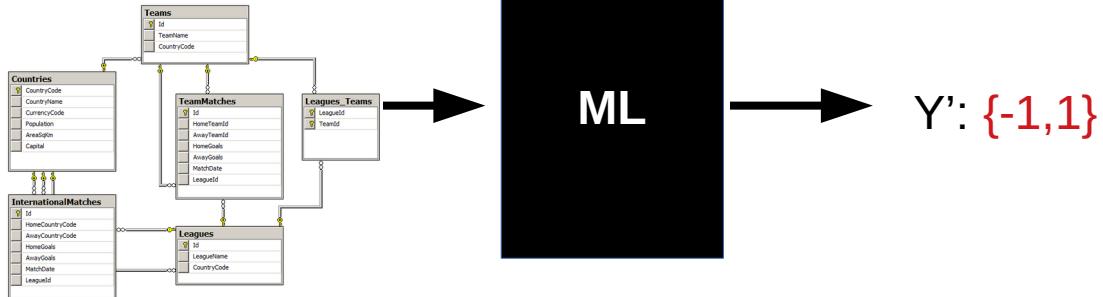
Where is the tensor data coming from?

- Problem data already comes as „clean“ vectors → „toy data“



Where is the tensor data coming from?

- Problem data already comes as „clean“ vectors → „toy data“
- Data comes from several relational data bases
 - select and merge data
 - need to find relevant variables
 - remove corrupt data
 - fill missing values
 - unify data formats
 - ...



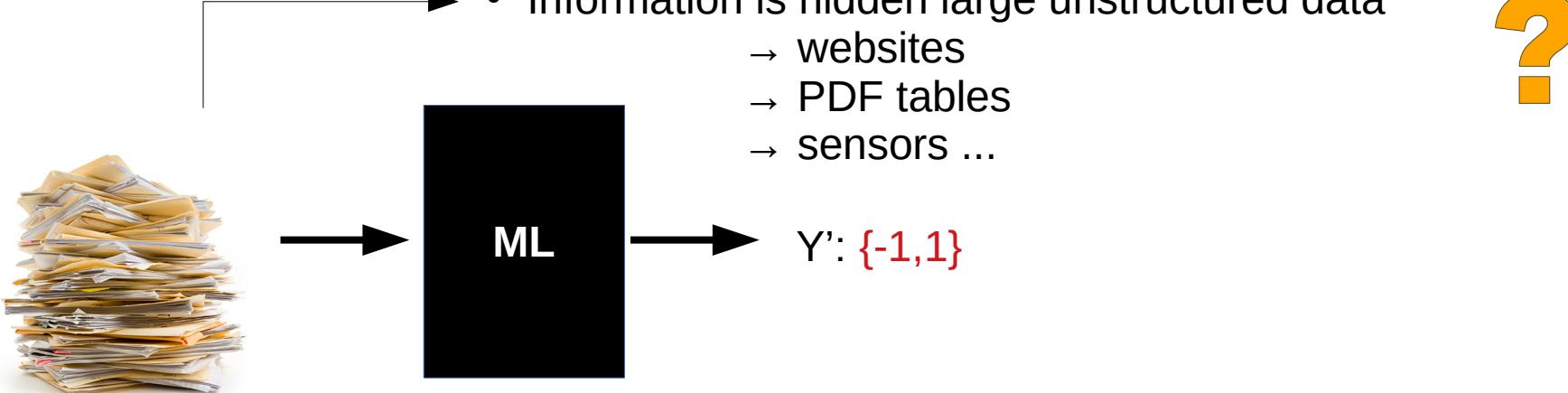
Where is the tensor data coming from?

- Problem data already comes as „clean“ vectors → „toy data“
- Data comes from several relational data bases
- Data is not numerical
 - categorical data, e.g. {„small“, „medium“, „large“}
 - text, images, audio
 - graphs
 - time series

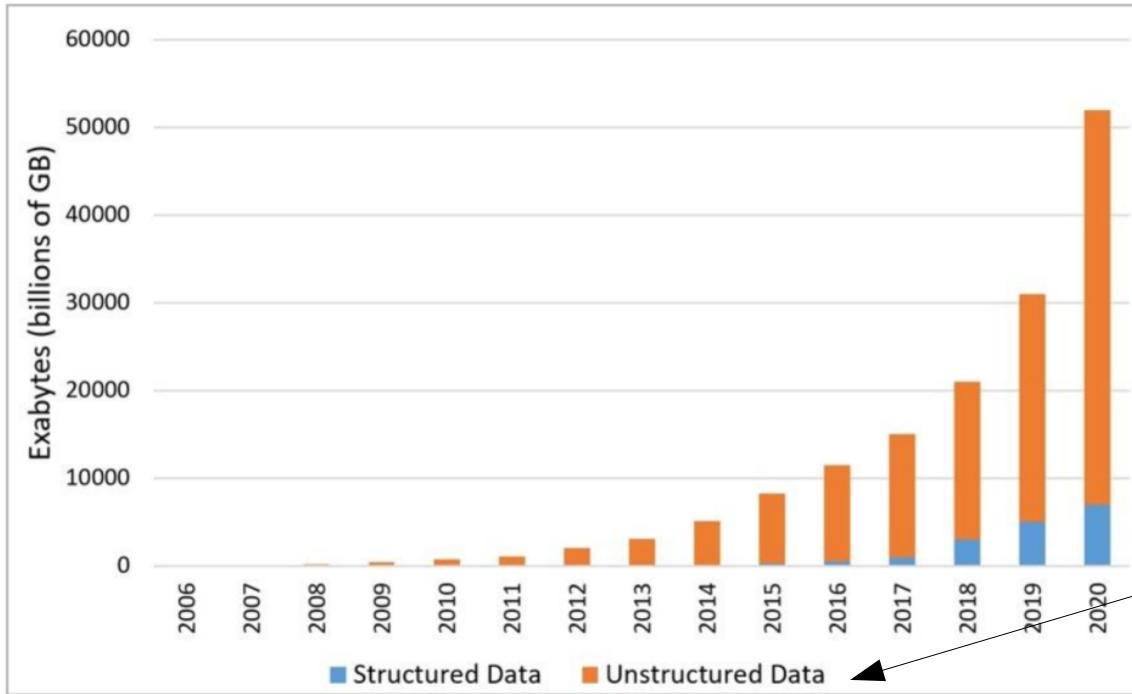


Where is the tensor data coming from?

- Problem data already comes as „clean“ vectors → „toy data“
- Data comes from several relational data bases
- Data is not numerical

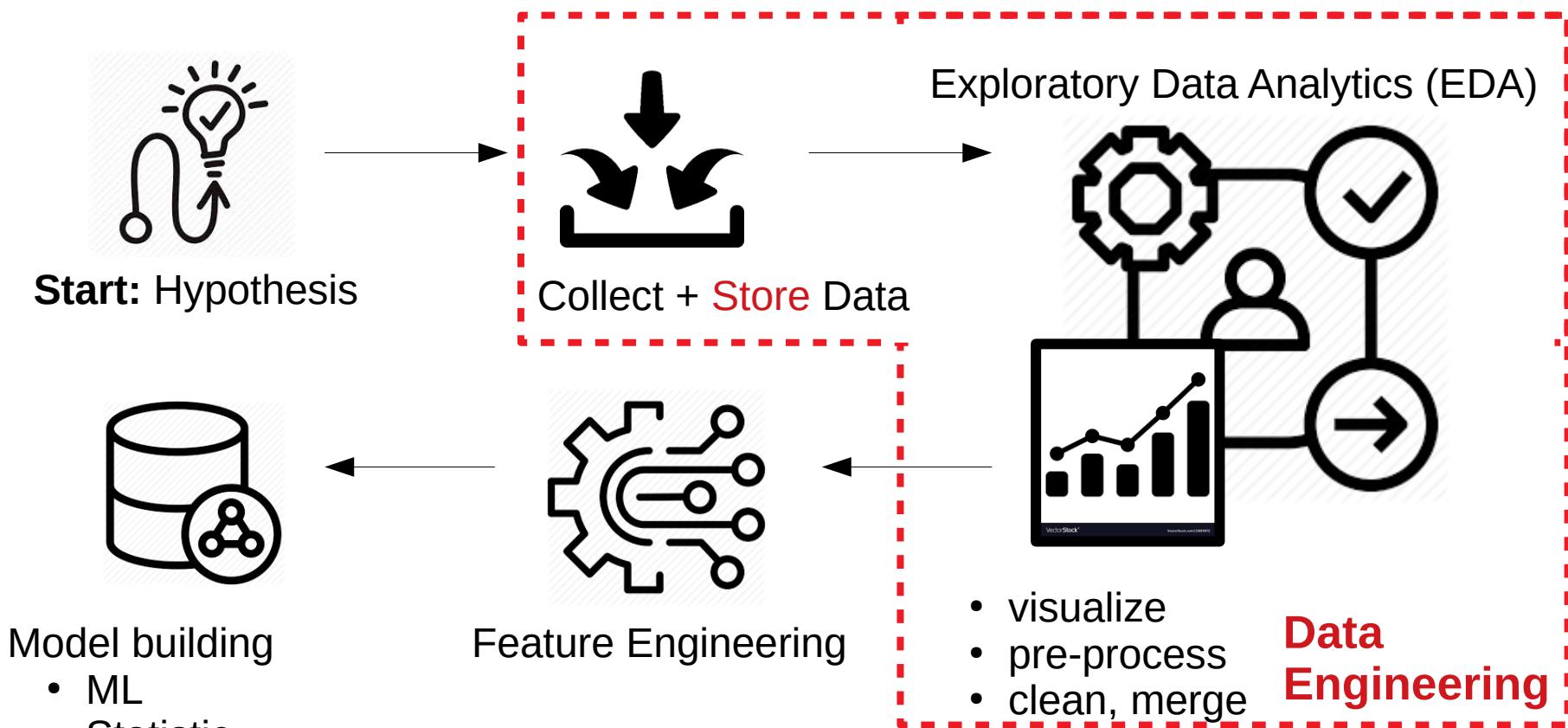


Where is the data coming from?



[3]

Not directly accessible by a
“query” and easily transformable
Into a tensor representation



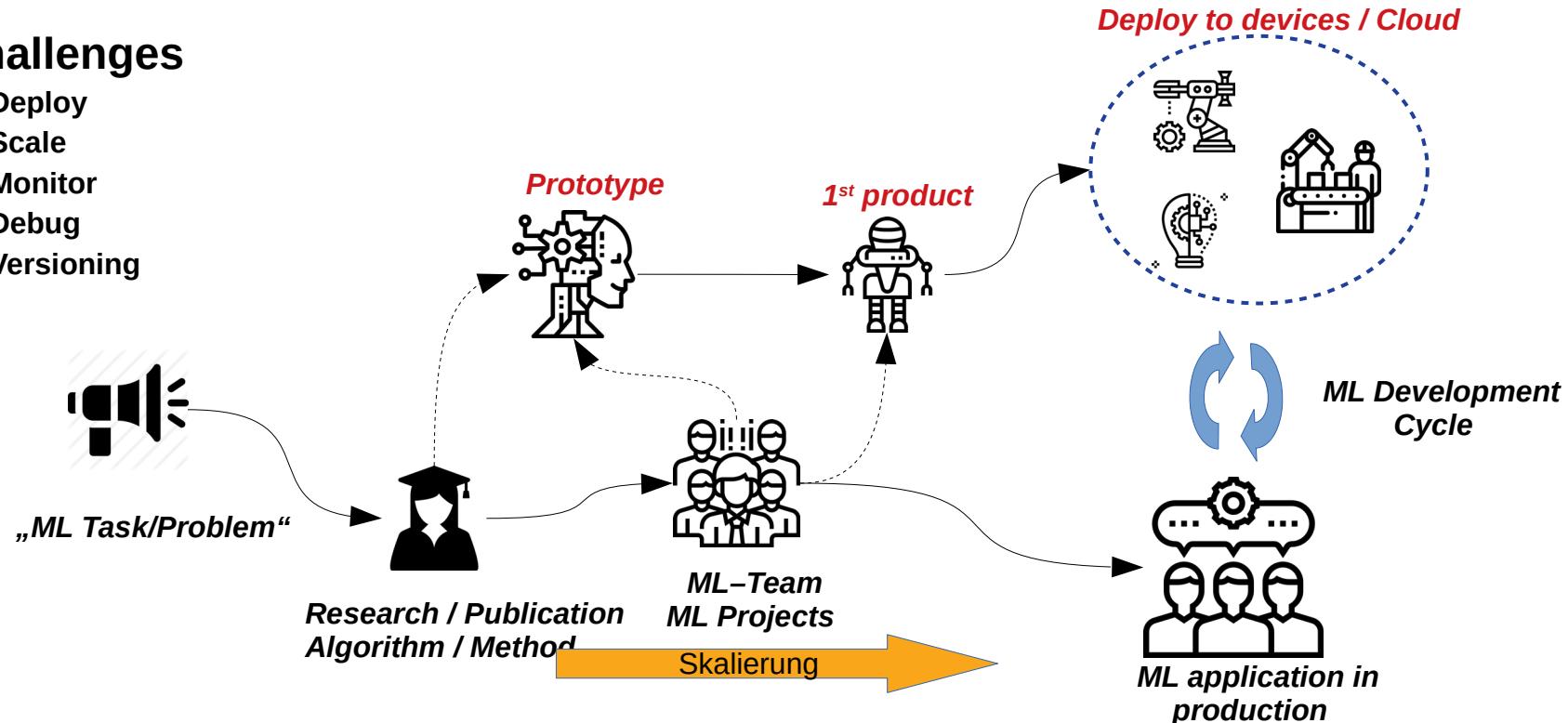
What is MLops?



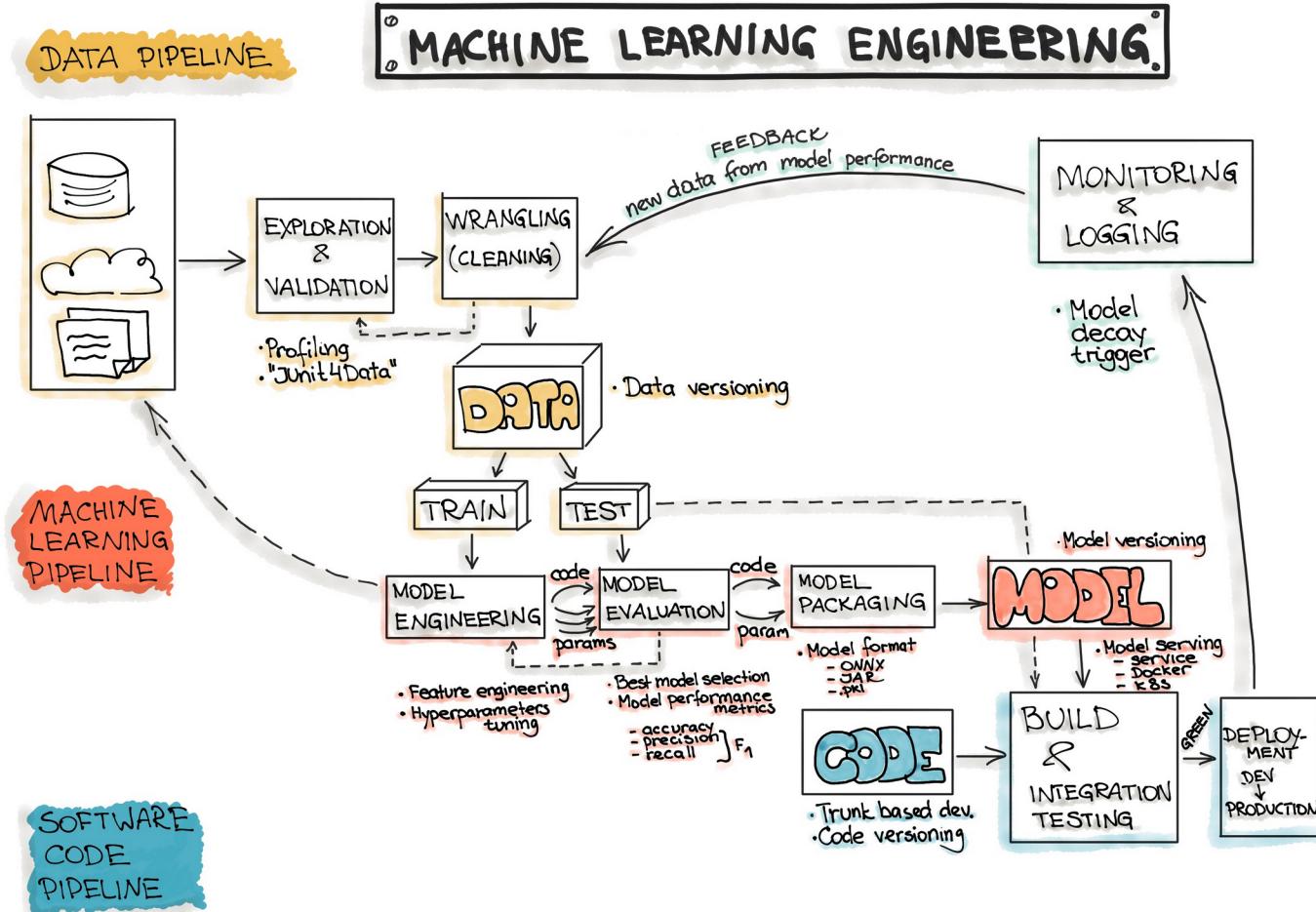
MLops → Machine Learning Operations !

Challenges

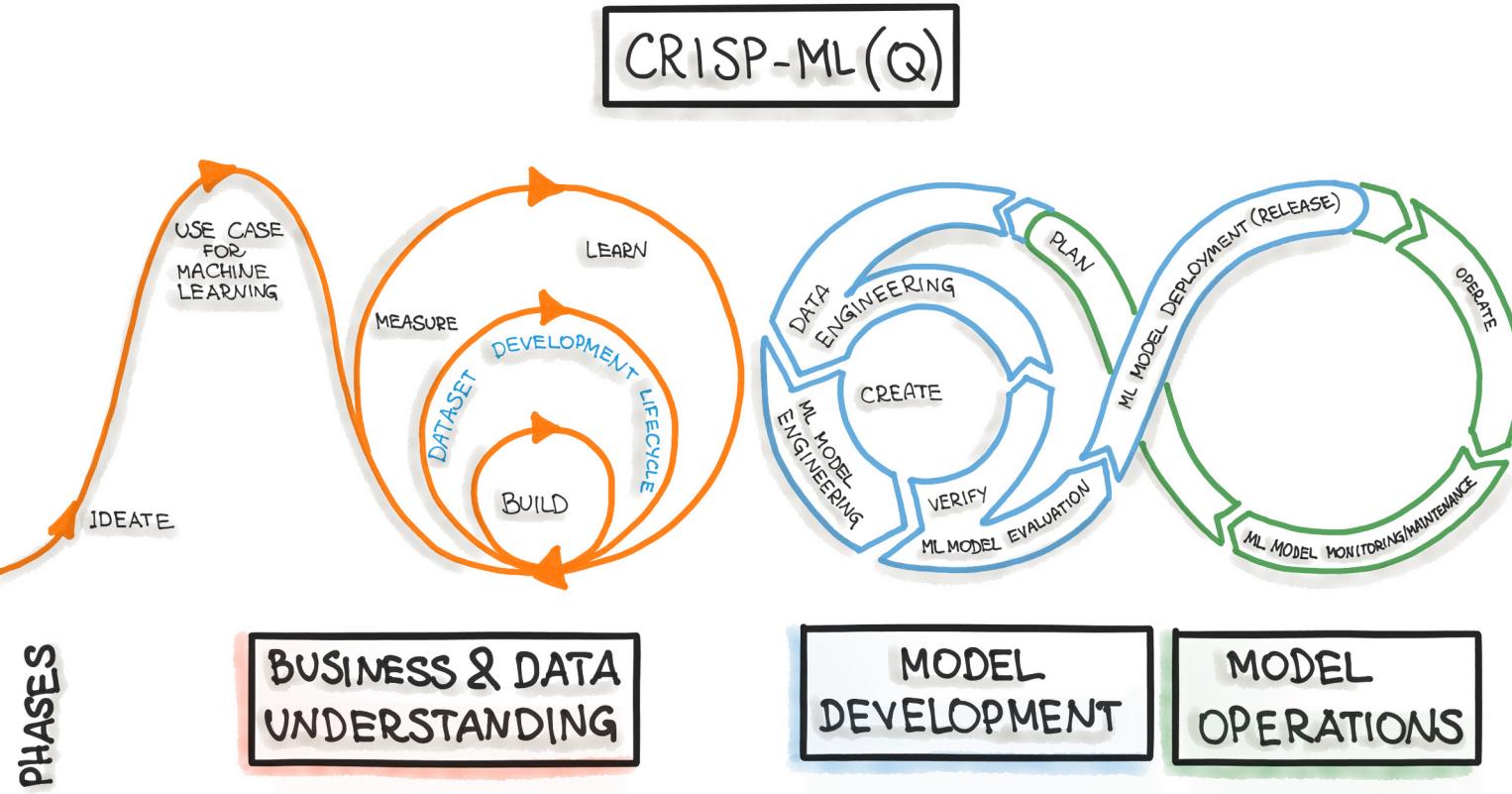
- Deploy
- Scale
- Monitor
- Debug
- Versioning



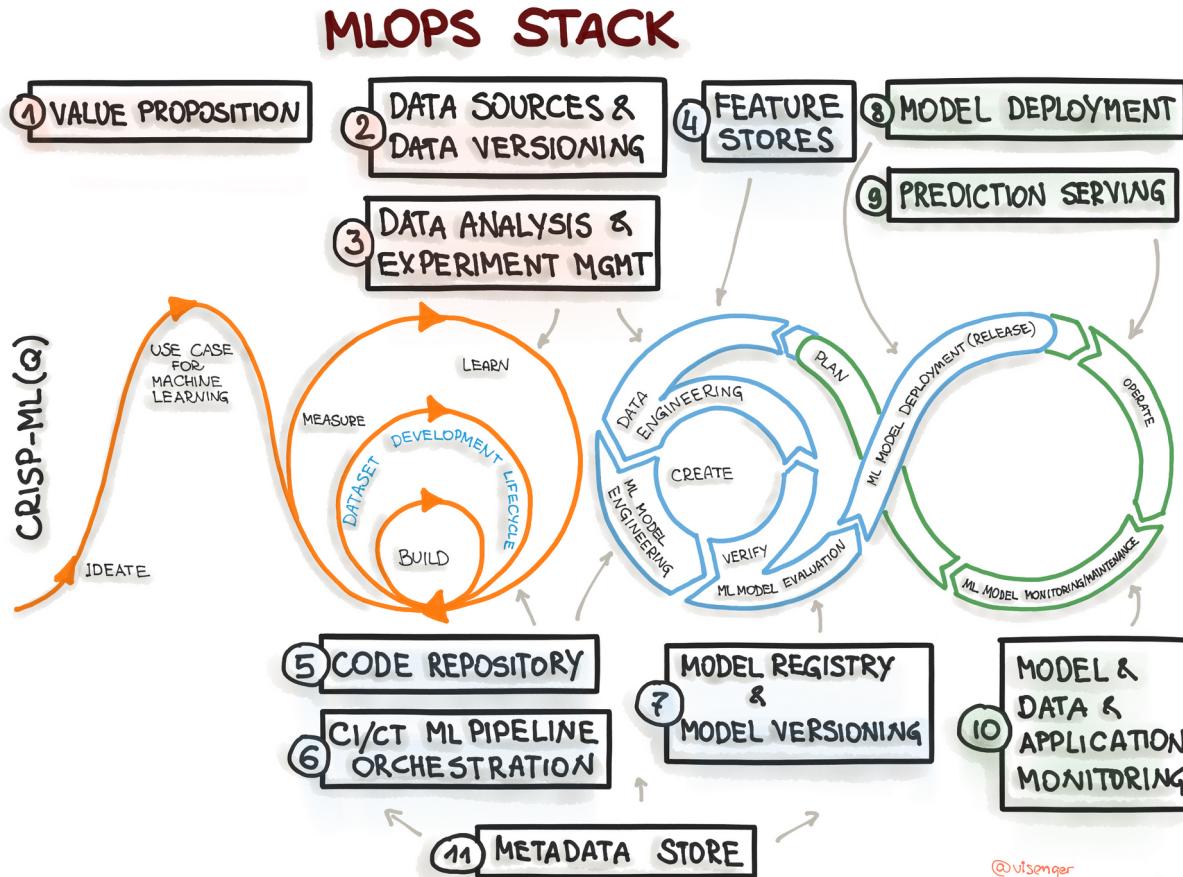
Machine Learning in REAL Applications



[<https://ml-ops.org/content/end-to-end-ml-workflow#data-engineering>]



@visenger



[<https://ml-ops.org/content/end-to-end-ml-workflow#data-engineering>]

Which programming Language did you use?

Training Hardware? → Cloud?

Which ML Frameworks?

Deployment Hardware?

How large was the Team?

Code Organization?

Model
Management?

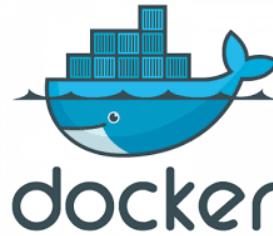
Which Scaling Framework?

Which NEW tools did you get to know?

Data Management?

MLOps

- Data management + versioning
- Experiment organization + tracking
- AutoML
- Model management, versioning and transfer
- Model deployment
- ML monitoring



docker



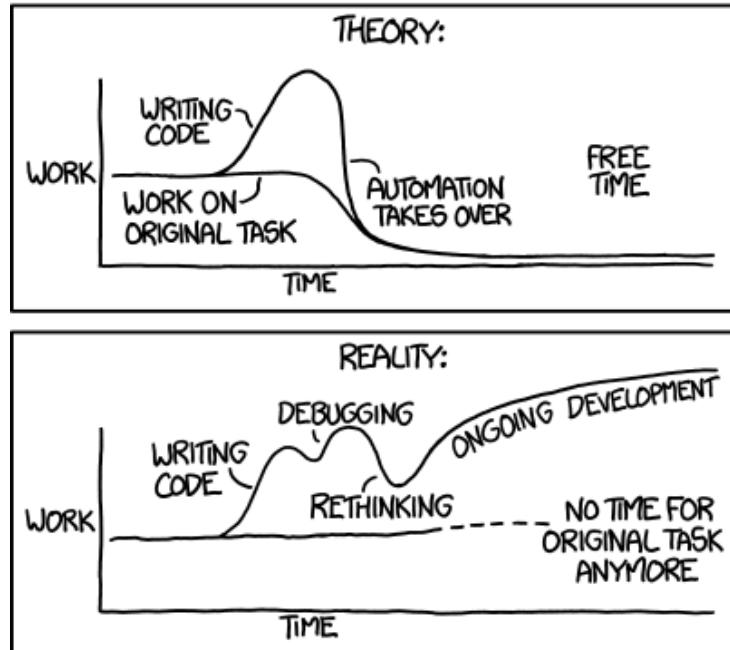
gradio



Deploying ML

- Containerization + Cloud
- Fast Web-APIs

"I SPEND A LOT OF TIME ON THIS TASK.
I SHOULD WRITE A PROGRAM AUTOMATING IT!"



[XKCD]

[1] free icons taken from <https://www.flaticon.com>

[2]

<https://www.forbes.com/sites/gilpress/2016/03/23/data-preparation-most-time-consuming-least-enjoyable-data-science-task-survey-says/>

[3] https://www.eetimes.com/author.asp?section_id=36&doc_id=1330462#