# Methodological topics
# Data-science specifics (part 1)
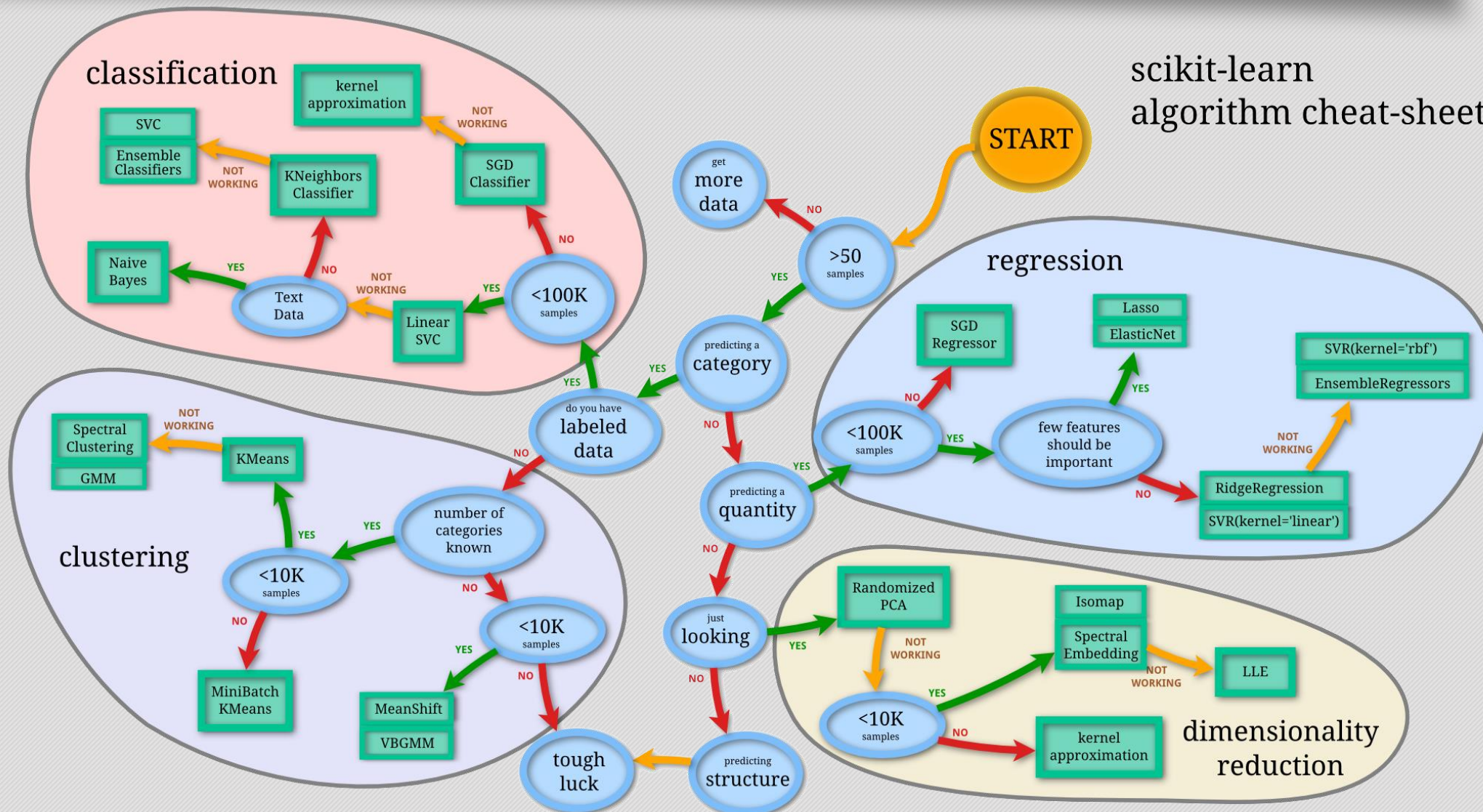
Ivan Srba

10th October 2018
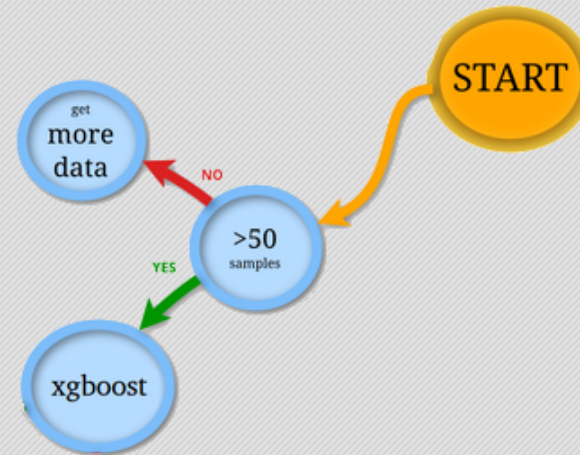
datalys

PeWe@FIIT
personalized web group

scikit-learn
algorithm cheat-sheet

**classification**

kernel approximation

SVC

Ensemble Classifiers

KNeighbors Classifier

SGD Classifier

NOT WORKING

Naive Bayes

Text Data

Linear SVC

<100K samples

NOT WORKING

NO

YES

YES

NOT WORKING

NO

START

get more data

NO

>50 samples

YES

predicting a category

do you have labeled data

YES

YES

NO

predicting a quantity

NO

just looking

NO

predicting structure

tough luck

**regression**

SGD Regressor

Lasso
ElasticNet

SVR(kernel='rbf')

EnsembleRegressors

NO

<100K samples

YES

few features should be important

YES

NOT WORKING

NO

RidgeRegression

SVR(kernel='linear')

**clustering**

Spectral Clustering

GMM

NOT WORKING

KMeans

number of categories known

YES

YES

<10K samples

NO

NO

<10K samples

YES

NO

MiniBatch KMeans

MeanShift

VBGMM

**dimensionality reduction**

Randomized PCA

NOT WORKING

Isomap

Spectral Embedding

NOT WORKING

LLE

YES

YES

<10K samples

NO

kernel approximation

scikit-learn
algorithm cheat-sheet

scikit-learn
algorithm cheat-sheet

http://scikit-learn.org/stable/tutorial/machine_learning_map/

# Warming-up

- Everything said last week applies perfectly also in case of all theses in data science domain

# Warming-up

- Everything said last week applies perfectly also in case of all theses in data science domain

- Summary of gold rules
  - Search for sources (research articles), organize them by dedicated tools

# Warming-up

- Everything said last week applies perfectly also in case of all theses in data science domain

- Summary of gold rules
  - Search for sources (research articles), organize them by dedicated tools
  - Analyze the existing solutions, write notes, compare them

# Warming-up

- Everything said last week applies perfectly also in case of all theses in data science domain

- Summary of gold rules
  - Search for sources (research articles), organize them by dedicated tools
  - Analyze the existing solutions, write notes, compare them
  - Select few most related articles, describe them in very details
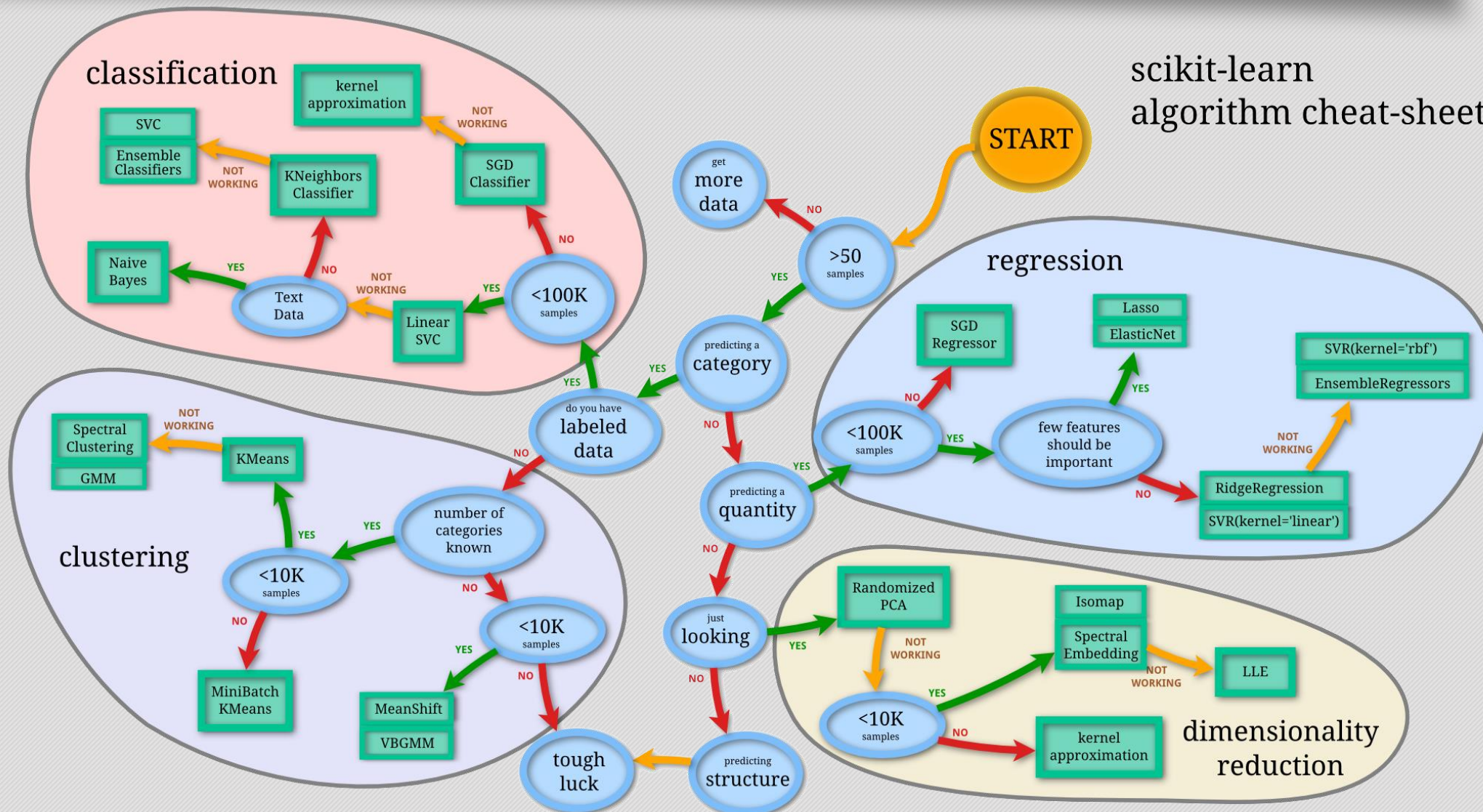
# Warming-up

- Everything said last week applies perfectly also in case of all theses in data science domain

- Summary of gold rules
  - Search for sources (research articles), organize them by dedicated tools
  - Analyze the existing solutions, write notes, compare them
  - Select few most related articles, describe them in very details
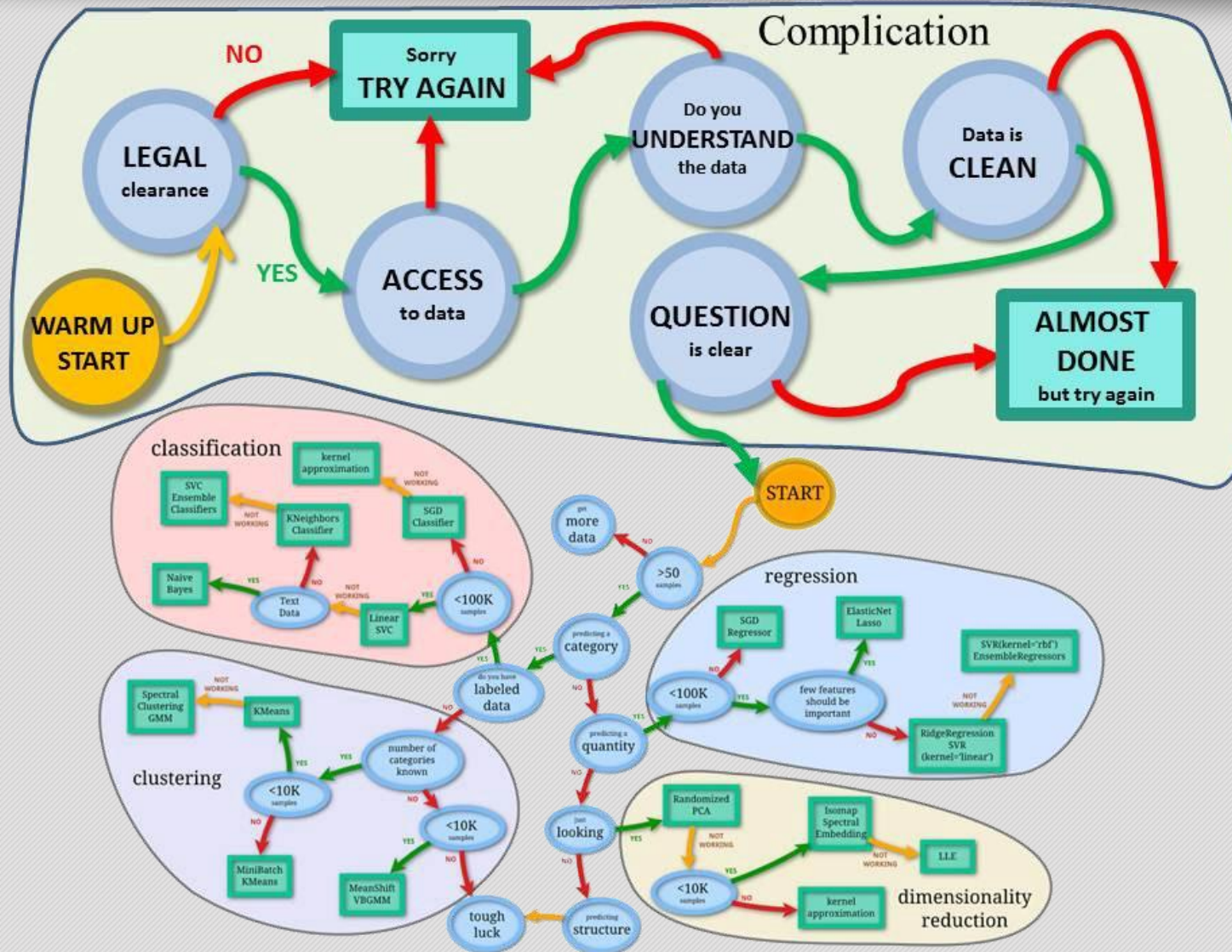  - Pay a strong attention to summary/discussion at the end of analyses' section

scikit-learn
algorithm cheat-sheet

http://scikit-learn.org/stable/tutorial/machine_learning_map/

https://medium.com/@chris_bour/an-extended-version-of-the-scikit-learn-cheat-sheet-5f46efc6cbb
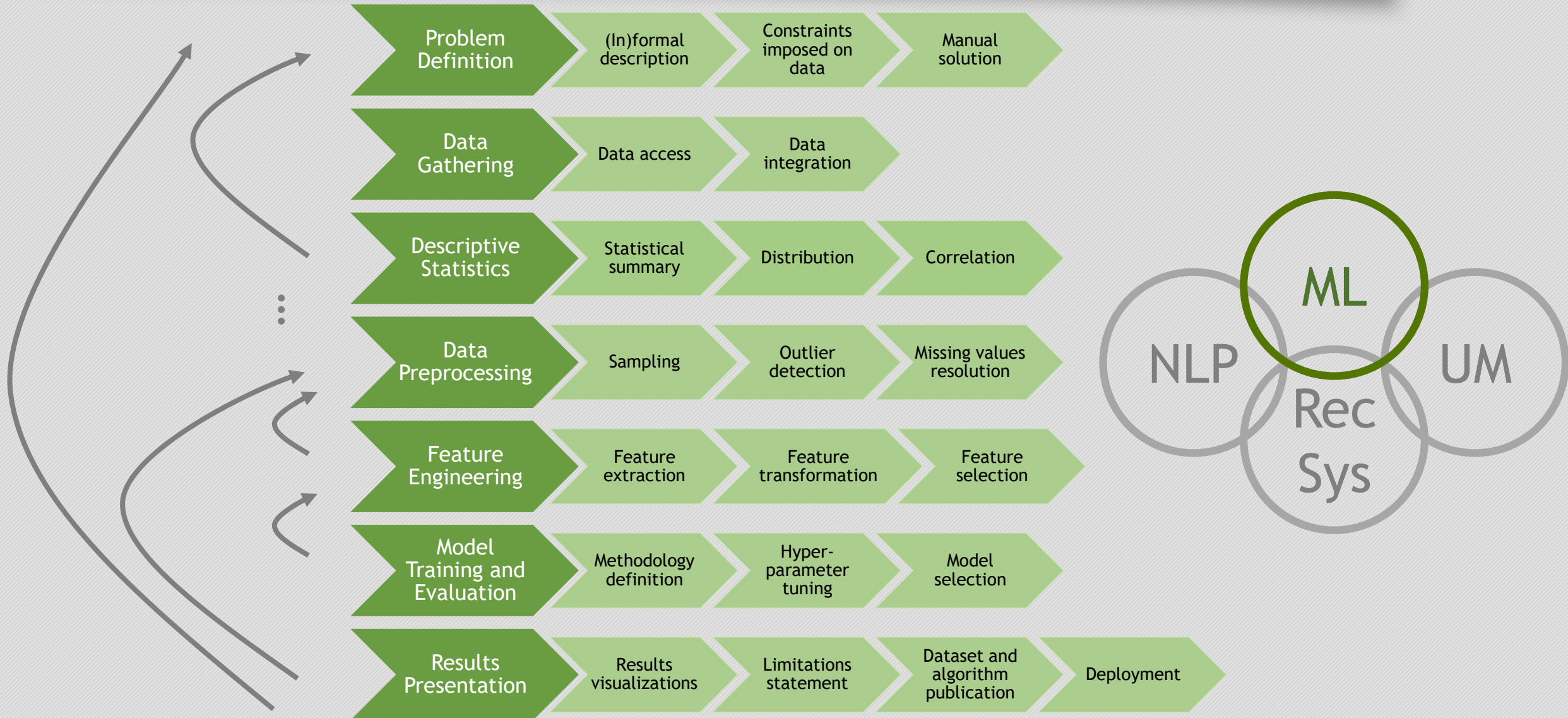
- … you need to answer before starting work on solution proposal and implementation:
    - How to define data-science (machine learning, …) task?
    - How to select/create appropriate dataset?
    - How to describe your dataset?
    - How to preprocess your dataset?
    - …

# Step 1: Problem Definition

15

- Clearly identify your problem you are trying to solve
  - Informal description
    - As you would explain it to your friends
    - Refer back to motivation stated in analyses' summary/discussion
  - Formal description
    - As a research question
    - As a hypothesis
    - As a machine learning task

| Problem Definition | (In)formal description | Constraints imposed on data | Manual solution |
|---|---|---|---|
| Data Gathering | Data access | Data integration | |
| Descriptive Statistics | Statistical summary | Distribution | Correlation |
| Data Preprocessing | Sampling | Outlier detection | Missing values resolution |
| Feature Engineering | Feature extraction | Feature transformation | Feature selection |
| Model Training and Evaluation | Methodology definition | Hyper-parameter tuning | Model selection |
| Results Presentation | Results visualizations | Limitations statement | Dataset and algorithm publication | Deployment |

- **Clearly identify your problem you are trying to solve**
  - Informal description
    - As you would explain it to your friends
    - Refer back to motivation stated in analyses' summary/discussion
  - Formal description
    - As a research question
    - As a hypothesis
    - As a machine learning task

- **Identify constraints imposed on required dataset**

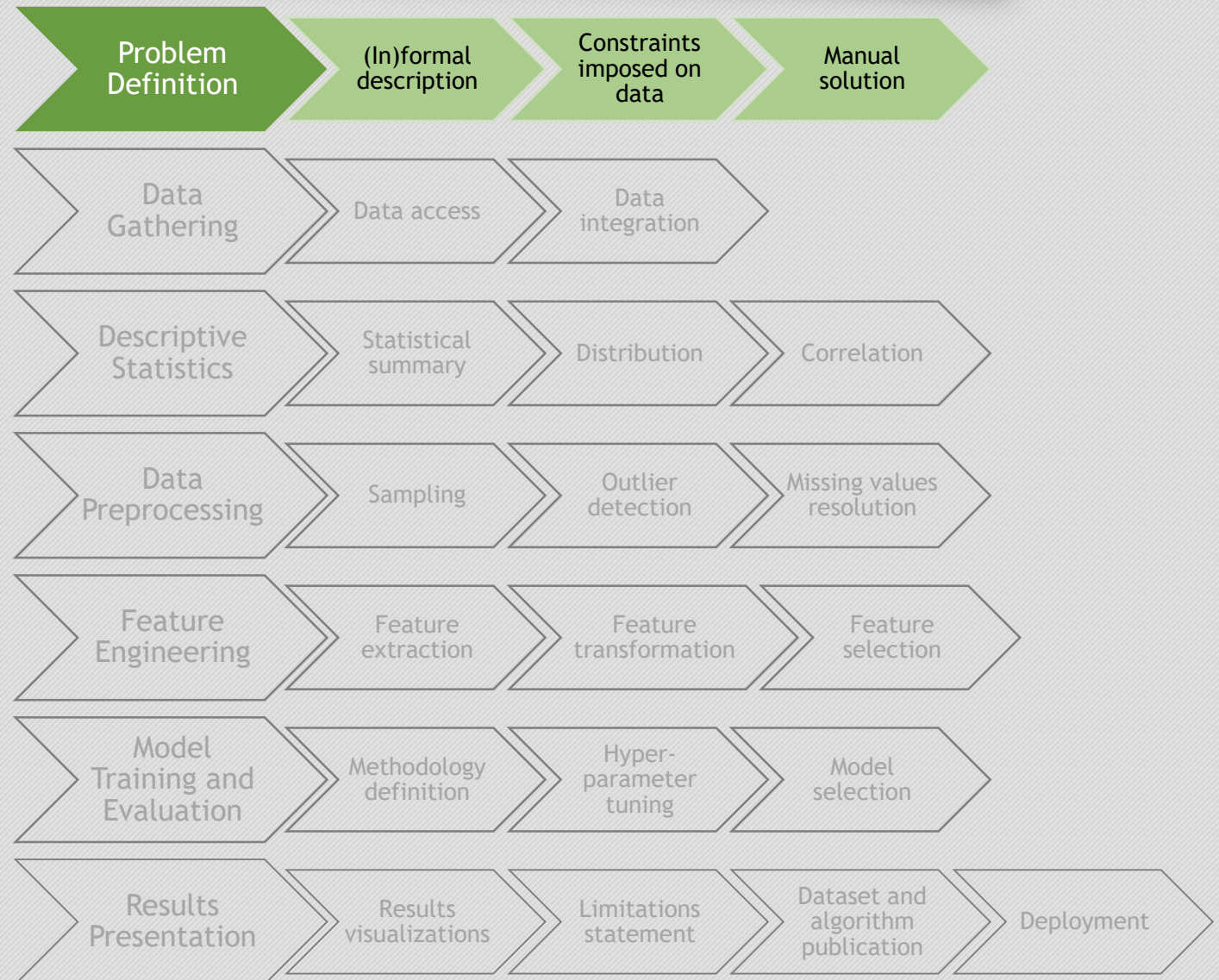| Problem Definition | (In)formal description | Constraints imposed on data | Manual solution |
|---|---|---|---|
| Data Gathering | Data access | Data integration | |
| Descriptive Statistics | Statistical summary | Distribution | Correlation |
| Data Preprocessing | Sampling | Outlier detection | Missing values resolution |
| Feature Engineering | Feature extraction | Feature transformation | Feature selection |
| Model Training and Evaluation | Methodology definition | Hyper-parameter tuning | Model selection |
| Results Presentation | Results visualizations | Limitations statement | Dataset and algorithm publication | Deployment |

- Clearly identify your problem you are trying to solve
  - Informal description
    - As you would explain it to your friends
    - Refer back to motivation stated in analyses' summary/discussion
  - Formal description
    - As a research question
    - As a hypothesis
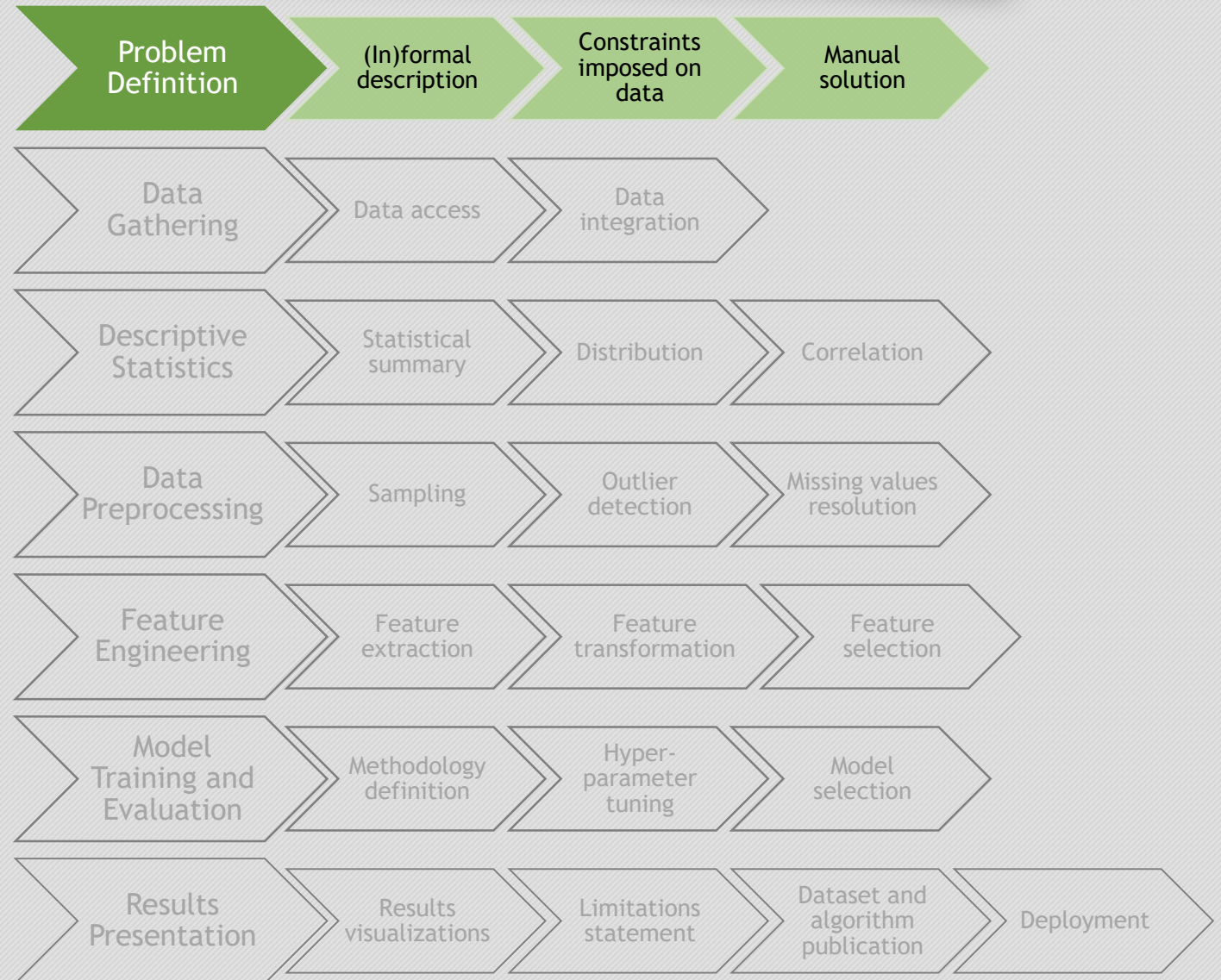    - As a machine learning task

- Identify constraints imposed on required dataset

- Explore possible manual solutions
  - If they do not exist, it is not a problem any more (in many cases)

| Problem Definition | (In)formal description | Constraints imposed on data | Manual solution |
|---|---|---|---|
| Data Gathering | Data access | Data integration | |
| Descriptive Statistics | Statistical summary | Distribution | Correlation |
| Data Preprocessing | Sampling | Outlier detection | Missing values resolution |
| Feature Engineering | Feature extraction | Feature transformation | Feature selection |
| Model Training and Evaluation | Methodology definition | Hyper-parameter tuning | Model selection |
| Results Presentation | Results visualizations | Limitations statement | Dataset and algorithm publication → Deployment |

- Data access
  - Prepared datasets, crawling, API
  - Legal issues

| Problem Definition | (In)formal description | Constraints imposed on data | Manual solution |
| --- | --- | --- | --- |

| Data Gathering | Data access | Data integration | |
| --- | --- | --- | --- |

| Descriptive Statistics | Statistical summary | Distribution | Correlation |
| --- | --- | --- | --- |
| Data Preprocessing | Sampling | Outlier detection | Missing values resolution |
| Feature Engineering | Feature extraction | Feature transformation | Feature selection |
| Model Training and Evaluation | Methodology definition | Hyper-parameter tuning | Model selection |
| Results Presentation | Results visualizations | Limitations statement | Dataset and algorithm publication / Deployment |

- **Data access**
  - Prepared datasets, crawling, API
  - Legal issues

- **Data integration**
  - Some ML task can be solved only when you integrate data from several sources
    - Different sources = different structure and format
  - Data consolidation
    - Entity mapping
    - User IDs (email, DB ID, cookie, username, …)
    - Item IDs (code, name, …)

| Problem Definition | (In)formal description | Constraints imposed on data | Manual solution |
|---|---|---|---|
| Data Gathering | Data access | Data integration | |
| Descriptive Statistics | Statistical summary | Distribution | Correlation |
| Data Preprocessing | Sampling | Outlier detection | Missing values resolution |
| Feature Engineering | Feature extraction | Feature transformation | Feature selection |
| Model Training and Evaluation | Methodology definition | Hyper-parameter tuning | Model selection |
| Results Presentation | Results visualizations | Limitations statement | Dataset and algorithm publication / Deployment |

- **Known your data!**
  - Otherwise, you are just guessing...

| Problem Definition | (In)formal description | Constraints imposed on data | Manual solution |
|---|---|---|---|

| Data Gathering | Data access | Data integration | |
|---|---|---|---|

| Descriptive Statistics | Statistical summary | Distribution | Correlation |
|---|---|---|---|

| Data Preprocessing | Sampling | Outlier detection | Missing values resolution |
|---|---|---|---|

| Feature Engineering | Feature extraction | Feature transformation | Feature selection |
|---|---|---|---|

| Model Training and Evaluation | Methodology definition | Hyper-parameter tuning | Model selection |
|---|---|---|---|

| Results Presentation | Results visualizations | Limitations statement | Dataset and algorithm publication | Deployment |
|---|---|---|---|---|

- **Known your data!**
  - Otherwise, you are just guessing…

- **Summarize data**
  - Volume of data (attributes, instances)
  - Data types
  - Distribution of data
  - Relations in data

- **Visualize data**
  - Histograms, boxplots, scatterplots

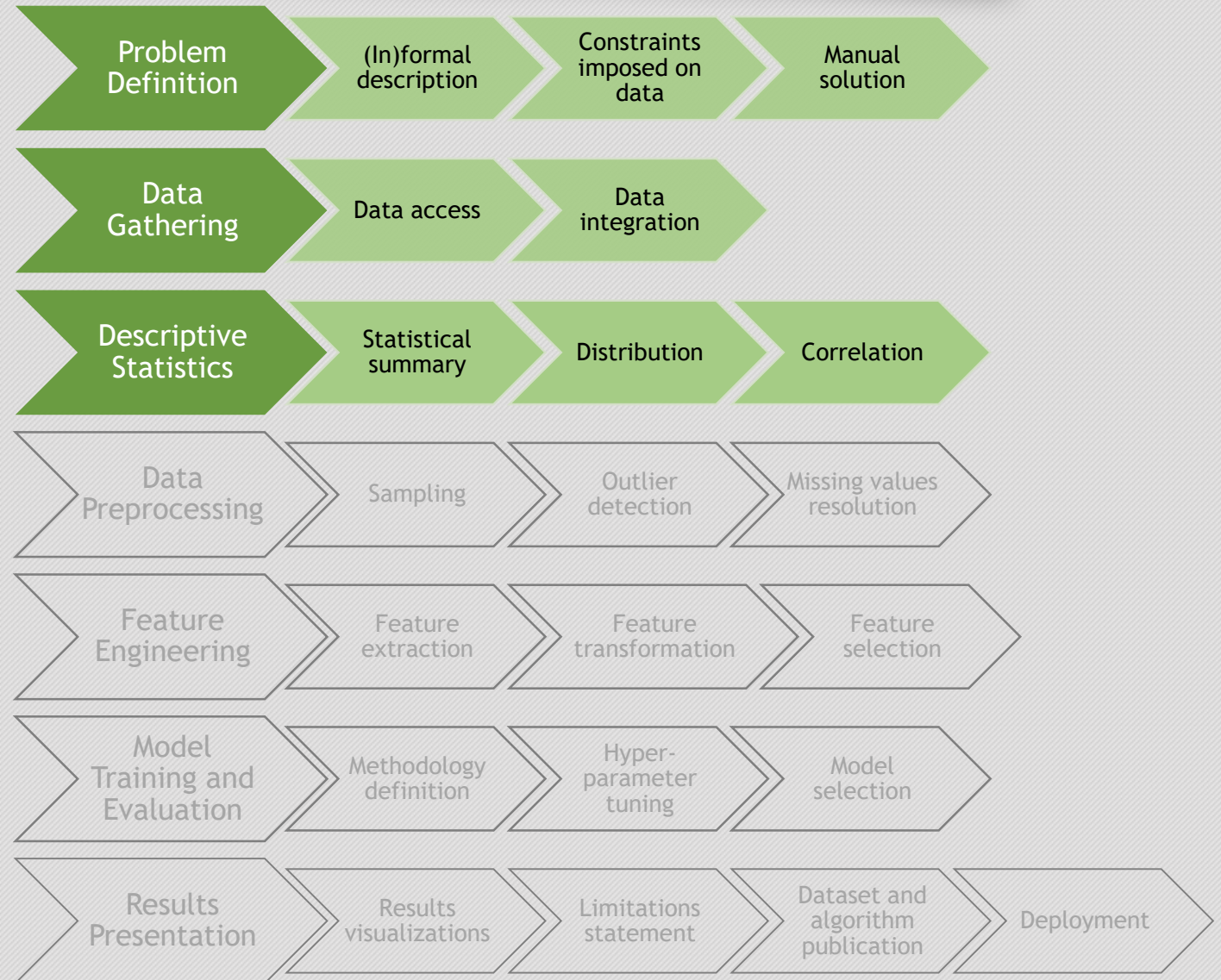| Problem Definition | (In)formal description | Constraints imposed on data | Manual solution |
|---|---|---|---|
| Data Gathering | Data access | Data integration | |
| Descriptive Statistics | Statistical summary | Distribution | Correlation |
| Data Preprocessing | Sampling | Outlier detection | Missing values resolution |
| Feature Engineering | Feature extraction | Feature transformation | Feature selection |
| Model Training and Evaluation | Methodology definition | Hyper-parameter tuning | Model selection |
| Results Presentation | Results visualizations | Limitations statement | Dataset and algorithm publication | Deployment |

- Known your data!
  - Otherwise, you are just guessing...

- Summarize data
  - Volume of data (attributes, instances)
  - Data types
  - Distribution of data
  - Relations in data

- Visualize data
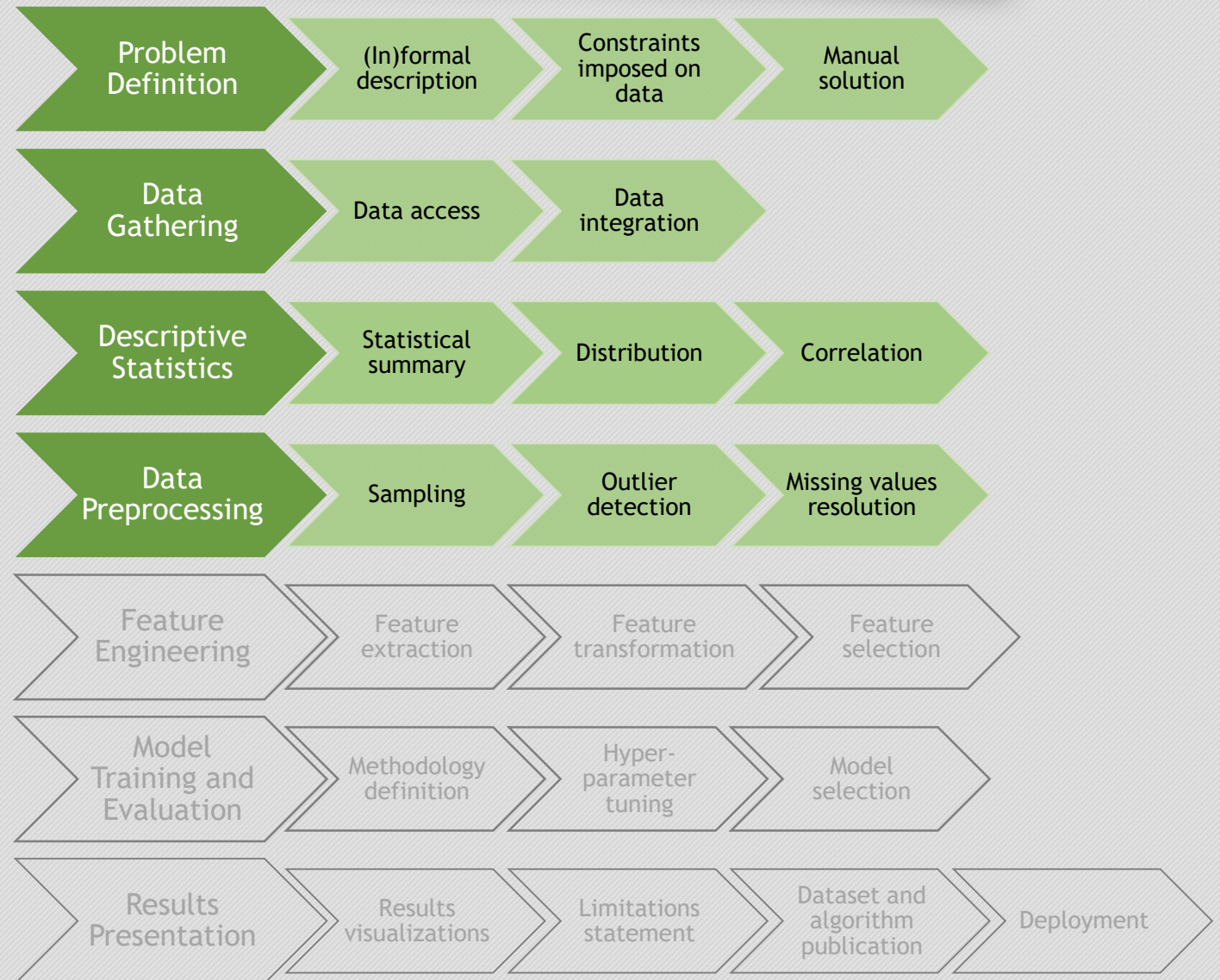  - Histograms, boxplots, scatterplots

- Result of descriptive statistics is an **important** input to all consequent steps

| Problem Definition | (In)formal description | Constraints imposed on data | Manual solution |
|---|---|---|---|
| Data Gathering | Data access | Data integration | |
| Descriptive Statistics | Statistical summary | Distribution | Correlation |
| Data Preprocessing | Sampling | Outlier detection | Missing values resolution |
| Feature Engineering | Feature extraction | Feature transformation | Feature selection |
| Model Training and Evaluation | Methodology definition | Hyper-parameter tuning | Model selection |
| Results Presentation | Results visualizations | Limitations statement | Dataset and algorithm publication |

Deployment

- Have large data? Do sampling!
  - Less data result in shorter training times
  - You can still finally run the model on larger portion of data

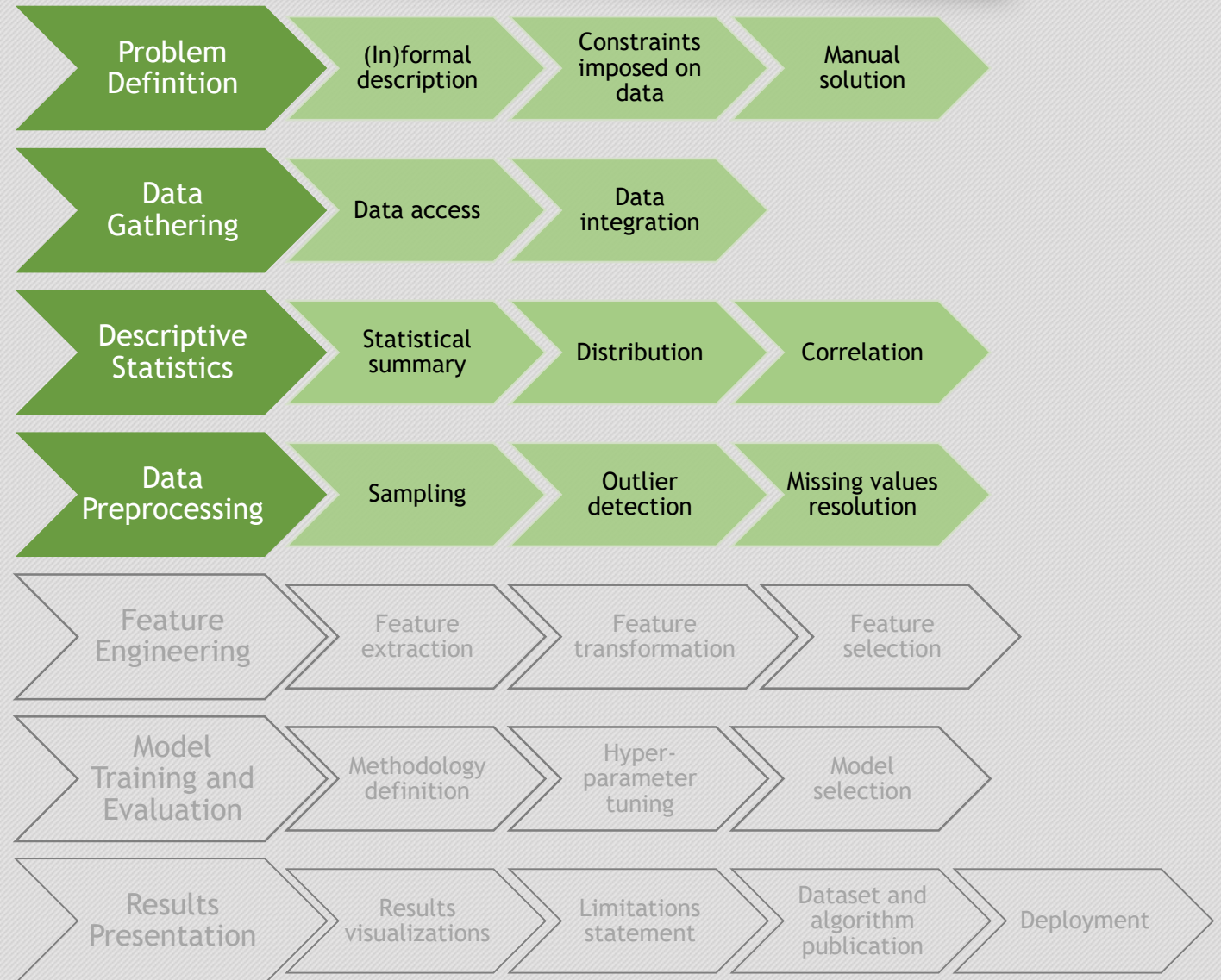| Problem Definition | (In)formal description | Constraints imposed on data | Manual solution |
| --- | --- | --- | --- |
| Data Gathering | Data access | Data integration | |
| Descriptive Statistics | Statistical summary | Distribution | Correlation |
| Data Preprocessing | Sampling | Outlier detection | Missing values resolution |
| Feature Engineering | Feature extraction | Feature transformation | Feature selection |
| Model Training and Evaluation | Methodology definition | Hyper-parameter tuning | Model selection |
| Results Presentation | Results visualizations | Limitations statement | Dataset and algorithm publication | Deployment |

- Have large data? Do sampling!
  - Less data result in shorter training times
  - You can still finally run the model on larger portion of data
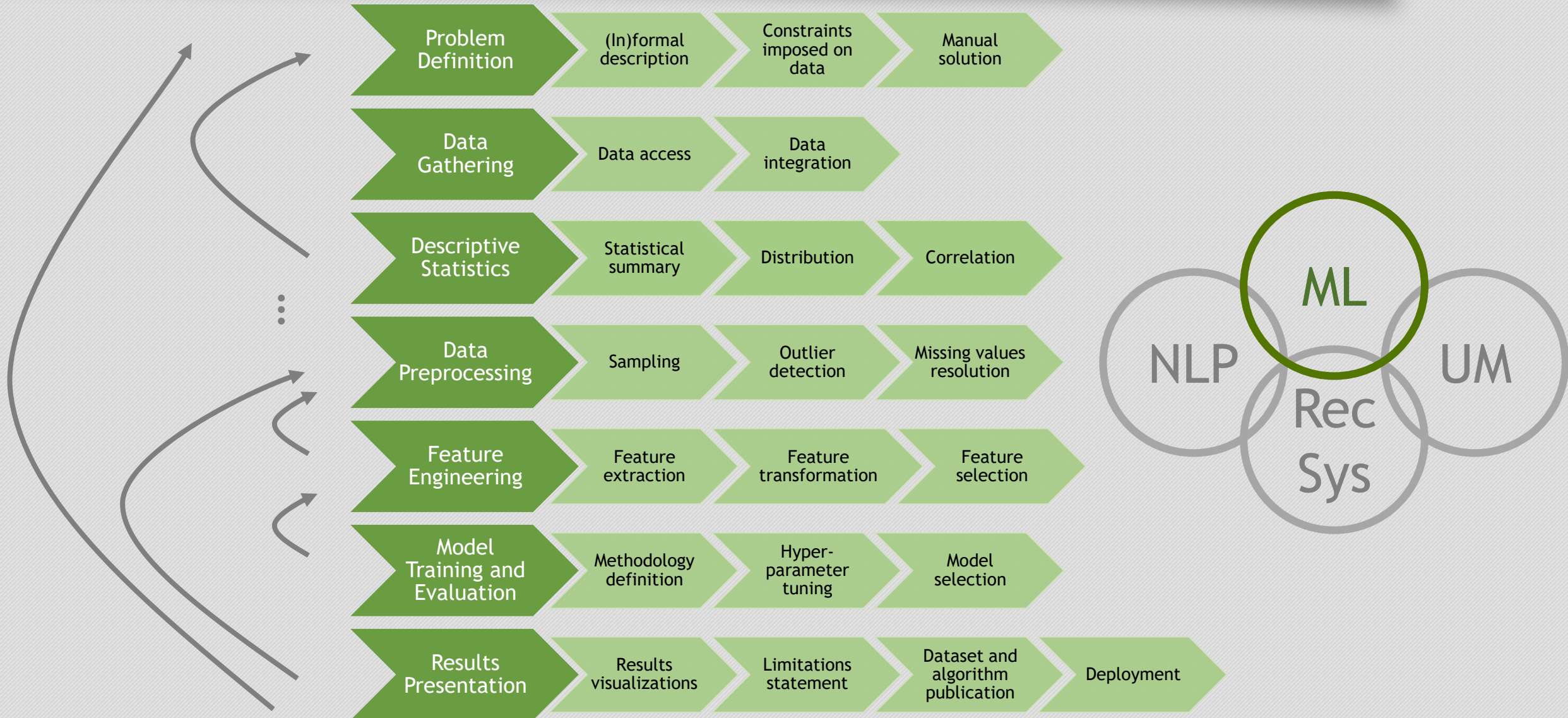
- Machine learning requires well-prepared data
  - Detect outliers
  - Replace missing values

| Problem Definition | (In)formal description | Constraints imposed on data | Manual solution |
|---|---|---|---|
| Data Gathering | Data access | Data integration | |
| Descriptive Statistics | Statistical summary | Distribution | Correlation |
| Data Preprocessing | Sampling | Outlier detection | Missing values resolution |
| Feature Engineering | Feature extraction | Feature transformation | Feature selection |
| Model Training and Evaluation | Methodology definition | Hyper-parameter tuning | Model selection |
| Results Presentation | Results visualizations | Limitations statement | Dataset and algorithm publication |

Deployment

- https://machinelearningmastery.com/4-steps-to-get-started-in-machine-learning/