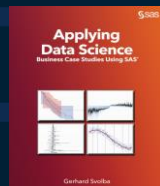


Hands-On Workshop

Machine Learning in Action: Feature Engineering, Modellierung und Modell Validierung in SAS Viya - Der Programmier-Ansatz



Gerhard Svolba, Data Scientist
SAS Austria



Data Scientist @SAS - [Medium](#) [LinkedIn](#) [Github](#) [SAS-Books](#) [SAS Articles](#)
[Youtube](#) [DataPreparation4DataScience](#) [Data Science Use Cases](#)

Copyright © SAS Institute Inc. All rights reserved.

Hinweis

- Ein Beitrag mit Beispiel Code zum Thema dieses Vortrags ist in Vorbereitung. Der Link findet sich dann in dieser Sammlung
- Data Science and Data Preparation Article Overview by Gerhard
 - <https://communities.sas.com/t5/SAS-Communities-Library/Data-Science-and-Data-Preparation-Article-Overview-by-Gerhard/ta-p/727875>
- Bei Fragen vorab kontaktieren sie gerne the Autor unter
 - Sastools.by.gerhard@gmx.net

Statistics, Machine & Deep Learning



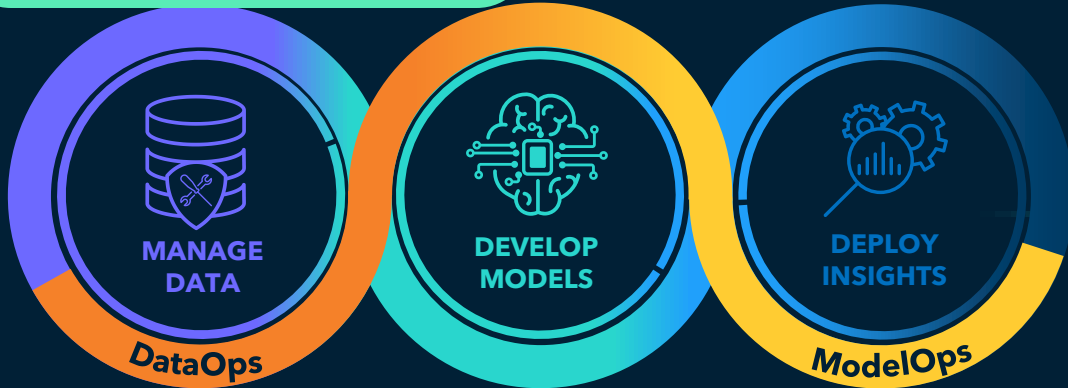
Natural Language Processing



Data Management



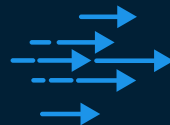
Deployment



Visualization



Decision Management



Computer & Machine Vision

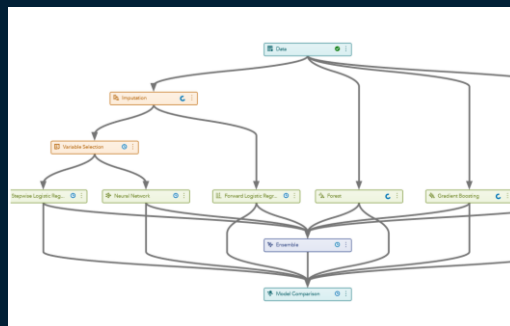
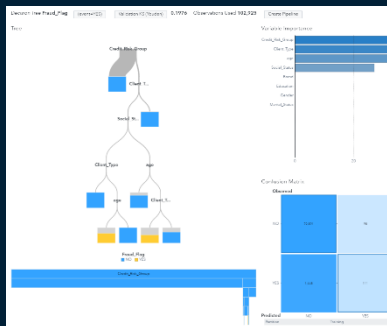


Forecasting, Optimization



Möglichkeiten der Interaktion mit der SAS Analytik Plattform

Graphische Benutzeroberfläche		Programmierung	
Visuelle Oberfläche	Model Studio	SAS	Open Source Sprache
Self-Service Analytik-Objekte Integration mit Model Studio & Model Manager	Pipelines und Knoten, Feature-Engineering, Optionen, Tuning, Open Source Integration, Integration mit dem Model Manager	Volle Flexibilität bei der Programmierung in der SAS Language (Procedures, Actions, Funktionen, ...) Open Source Integration	Interaktion mit der SAS Analytik- Plattform aus dem Jupyter- Notebook oder R-Studio heraus



```

28 proc gradboost data=cas1.fc_review
29   earlystop(tolerance=0 stagnation=5)
30   numbin=20 binmethod=BUCKET
31   maxdepth=5
32   maxbranch=2
33   minleafsize=5
34   assignmissing=USEINSEARCH minuseinsearch=1
35   seed=12345
36   printtarget;
37 ;
38 ;
39 partition relevarv=partind. (TRAIN='1' VALIDATE='0');
40 autotune useparameters=CUSTOM tuningparameters=(
41   lasso(LB=0 UB=10 INIT=0)
42   learningrate(LB=0.01 UB=1 INIT=0.1)
43   ntrees(LB=20 UB=150 INIT=100)
44   ridge(LB=0 UB=10 INIT=0)
45   samplingrate(LB=0.1 UB=1 INIT=0.5)
46   vars_to_try(LB=1 UB=7 INIT=7)
47 )
48 searchmethod=GA objective=KS maxtime=500
49 maxevals=50 maxiters=5 popsize=10
50 targetevents='1'
51 ;
    
```

```

from swei_loader import render_html
from swei import *
from pprint import pprint

import matplotlib.pyplot as plt
import pandas as pd
from pandas import *
import numpy as np

import seaborn as sns
get_ipython().magic('matplotlib inline')
from IPython.core.interactiveshell import InteractiveShell
InteractiveShell.ast_node_interactivity = "all"

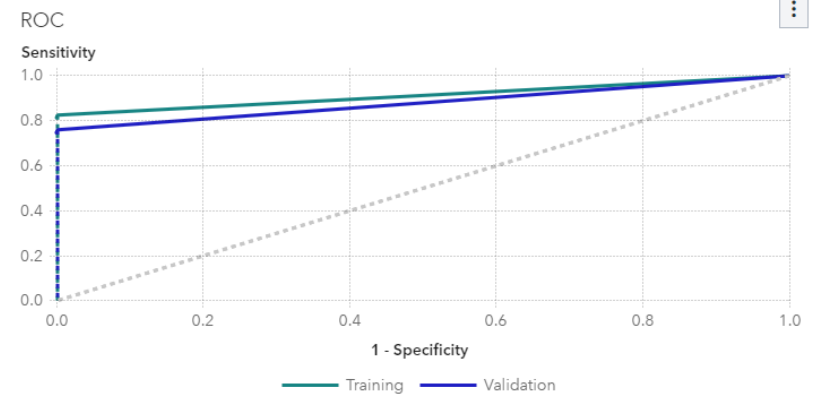
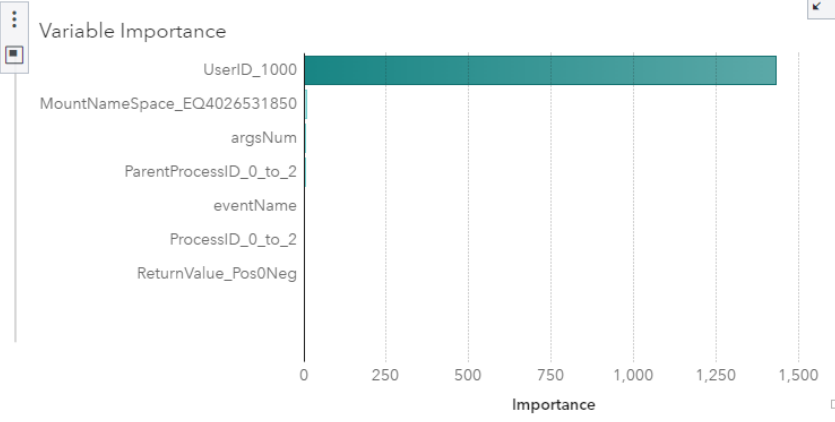
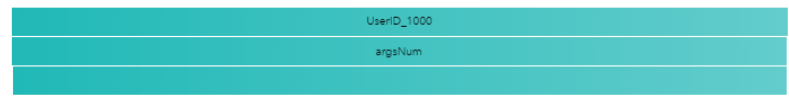
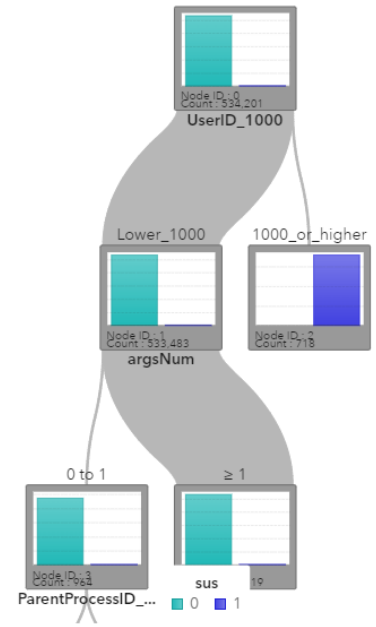
# Provide Connection Information and Updated Data if not yet available
casHost = 'cas1-cas01-01.sas.com'
casPort = 8779
casAuth = 'sas7catcat'
casDataDir = '/sas7catcat'
casTable = 'sas7catcat'

# Create Demo Instance SAS Club and Load Actions Sets for Decision Trees
SASClub = SASClub(casHost, casPort, casAuth, casDataDir, casTable)
SASClub.loadactionset(actionset='DecisionTree')
if not SASClub.table.exists(table='indata').exists():
    'tbl' = SASClub.upload_file(indata_dir='/indata', casout=('name', indata))
NOTE: Added action set 'DecisionTree'.
    
```

Decision Tree sus (event=1) Validation C Statistic 0.880 Observations Used 763,144

Create Pipeline

Tree



Nodes

Filter

Data Mining Preprocessing

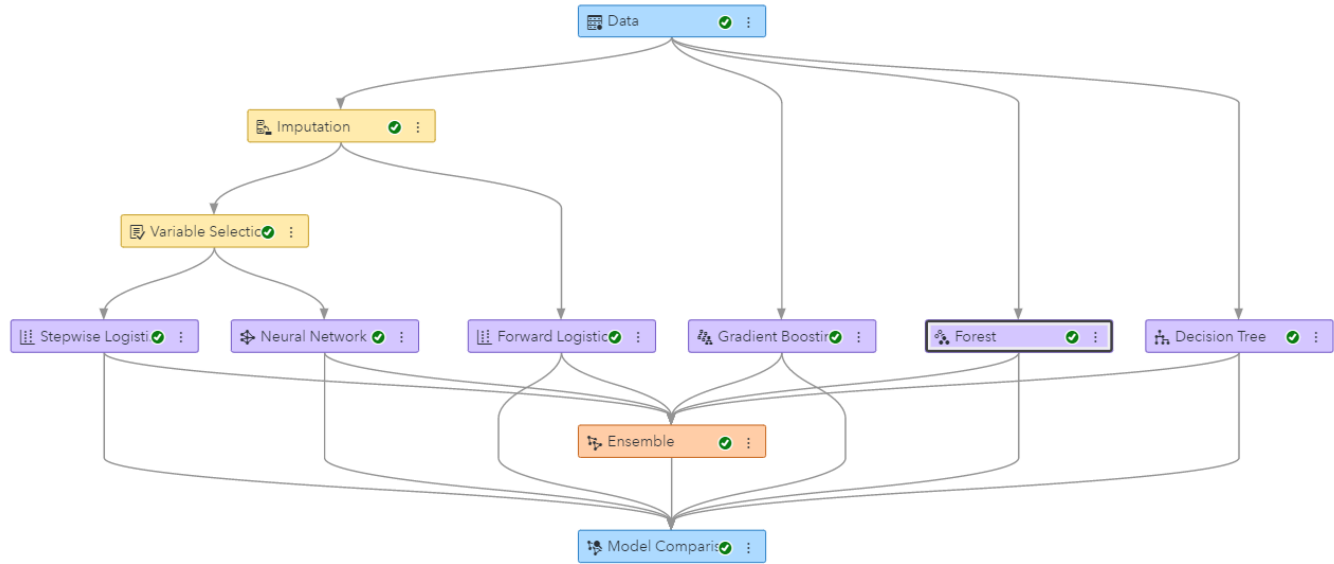
- Anomaly Detection
- Clustering
- Feature Extraction
- Feature Machine
- Filtering
- Imputation
- Interactive Grouping
- Manage Variables
- PCA1
- Reject Inference
- Replacement
- Text Mining
- Transformations
- Variable Clustering
- Variable Selection

Supervised Learning

- Batch Code
- Bayesian Network
- Decision Tree
- Forest
- Forest_1
- Forest_2
- Forest_3

SVM from VisualAnalytics AdvTemplate

Run Pipeline



Context: Cyber Security

- ... anomalous data points and shifts in the data distribution are inevitable. From a cyber security perspective, these anomalies and **dataset shifts are driven by both defensive and adversarial advancement**.
- To withstand the cost of critical system failure, the **development of robust models is therefore key to the performance, protection, and longevity of deployed defensive systems**.
- We present the **BPF-extended tracking honeypot (BETH) dataset** as the first cybersecurity dataset for uncertainty and robustness benchmarking. Collected using a novel honeypot tracking system, our dataset has the following properties that make it attractive for the development of robust ML methods:
 1. At over eight million data points, this is one of the largest cyber security datasets available
 2. It contains modern host activity and attacks
 3. It is fully labelled
 4. It contains highly structured but heterogeneous features
 5. Each host contains benign activity and at most a single attack, which is ideal for behavioural analysis and other research tasks. In addition to the described dataset

BETH Dataset: Real Cybersecurity Data for Anomaly Detection Research

FEATURE	TYPE	DESCRIPTION
TIMESTAMP	FLOAT	SECONDS SINCE SYSTEM BOOT
PROCESSID*	INT	INTEGER LABEL FOR THE PROCESS SPAWNING THIS LOG
THREADID	INT	INTEGER LABEL FOR THE THREAD SPAWNING THIS LOG
PARENTPROCESSID*	INT	PARENT'S INTEGER LABEL FOR THE PROCESS SPAWNING THIS LOG
USERID*	INT	LOGIN INTEGER ID OF USER SPAWNING THIS LOG
MOUNTNAMESPACE*	INT (LONG)	SET MOUNTING RESTRICTIONS THIS PROCESS LOG WORKS WITHIN
PROCESSNAME	STRING	STRING COMMAND EXECUTED
HOSTNAME	STRING	NAME OF HOST SERVER
EVENTID*	INT	ID FOR THE EVENT GENERATING THIS LOG
EVENTNAME	STRING	NAME OF THE EVENT GENERATING THIS LOG
ARGSNUM*	INT	LENGTH OF ARGS
RETURNVALUE*	INT	VALUE RETURNED FROM THIS EVENT LOG (USUALLY 0)
STACKADDRESSES	LIST OF INT	MEMORY VALUES RELEVANT TO THE PROCESS
ARGS	LIST OF DICTIONARIES	LIST OF ARGUMENTS PASSED TO THIS PROCESS
SUS	INT (0 OR 1)	BINARY LABEL AS A SUSPICIOUS EVENT (1 IS SUSPICIOUS, 0 IS NOT)
EVIL	INT (0 OR 1)	BINARY AS A KNOWN MALICIOUS EVENT (0 IS BENIGN, 1 IS NOT)

<https://www.kaggle.com/datasets/katehighnam/beth-dataset>



Pre-Processing / Feature Engineering

See Appendix A in the paper

processId: Process IDs 0, 1, and 2 are meaningful since these are always values used by the OS, but otherwise a random number is assigned to the process upon creation. We recommend replacing `processId` with a binary variable indicating whether or not `processId` is 0, 1, or 2.

threadId: While this value did not appear useful in our analysis, it might suggest how to link process calls if obfuscated in the system. No conversion is recommended at this time.

parentProcessId: Same as `processId`, the same mapping to a binary variable should suffice.

userId: The default in Linux systems is to assign OS activity to some number below 1000 (typically 0). As users login, they are assigned IDs starting at 1000, incrementally. This can be altered by a user, but none of the current logs gave evidence an attacker did this. We used a binary variable to indicate `userId < 1000` or `userId ≥ 1000`. Alternatively, one could use an ordinal mapping that buckets all `userId < 1000` at zero and then increment upwards for each new user. Also, no more than four logins were viewed per host in our current datasets.

BETH Dataset: Real Cybersecurity

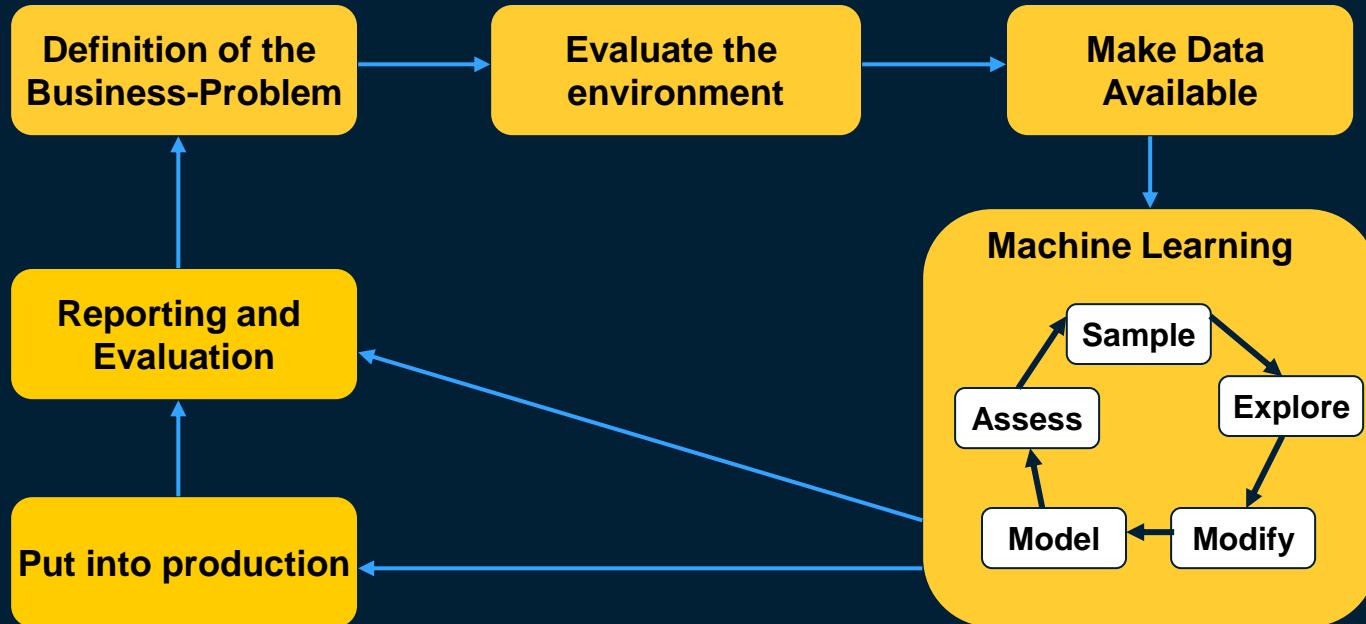
argsNum: This raw feature is included as-is, since, at this time, adequately parsing `args` requires either more sophisticated pre-processing or a more complex ML model.

returnValue: This is also called the exit status and can be used to determine whether a call completed successfully or not. Mappings for this can vary, as this value is decided between the parent and child process. We mapped `returnValue` into three values based on the common usage of the field: -1 when negative (bad errors), 0 when zero (success), and 1 when positive (success and signalling something to the parent process).

~~**stackAddresses:**~~ It is difficult to clearly relate this feature during manual analysis and the large values within a variable size list make processing automatically difficult without encoding or learning an extra embedding. Thus this field was dropped from training our baselines.

~~**args:**~~ There are many options in this variable list of dictionaries. For simplicity, we refrain from utilising any of these values. However, more features can and should be created for future work.

Machine Learning as part of a „Closed Loop“ Environment



Flow SAS Studio (Viya4)

The screenshot displays the SAS Studio interface for developing code and flows. The top navigation bar includes a menu icon, the title "SAS® Studio - Develop Code and Flows", and search, notification, and help icons. Below the navigation bar is a menu with "New", "Options", "View", "Open", and "Save All". The left sidebar contains an "Explorer" panel with a tree view of folders and files, including "My Favorites", "Folder Shortcuts", "My Snippets", "My Tasks", "SAS Videos", "Claims_Prediction", "share-data", "SAS Content", "AloT Projects", and "Conversational Flows". The main workspace shows a flow diagram titled "* CAS ML Programming HandsOn 1.2.flw". The flow consists of six steps: "1 Check and Sample Data", "2 Create Features", "3 Build Models", "4 Assess Models", "5 Score Fresh Data", and "6 Validate and Operationalize". The interface also includes a toolbar with "Run", "Cancel", and other action buttons, and a status bar at the bottom right showing the date and time: "Mar 30, 2023, 10:30:31 AM".

Flow SAS Studio (Viya4)

1 Check and Sample Data

2 Create Features



04 Bin
"Eventname"



Features
based on...



05 Variable
Importance



06 Screen
Variables



07 Feature
Machine

3 Build Models



11 Logistic
Regression



12 Logistic
Regression2...



12 Gradient
Boosting

4 Assess Models



21 PCreate
Assessment...



22 Create Lift
Charts

Ausgewählte SAS Procedures in SAS Viya

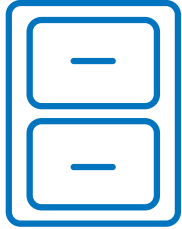
- [SAS Visual Statistics](#)

The screenshot shows the SAS Help Center interface. The top navigation bar includes 'SAS Help Center' and 'SAS Viya Platform Programming Documentation | 2023.03'. The left sidebar contains a navigation menu with categories like 'Welcome to SAS Programming Documentation', 'What's New', 'Learning SAS Viya Platform Programming', 'Syntax Quick Links', and 'SAS Procedures by Name and Product'. The 'SAS Procedures by Name and Product' section is expanded, showing a list of procedures: ASSESS, BART, BINNING, CARDINALITY, CORRELATION FREQTAB, GAMMOD, GAMSELECT, GENSELECT, ICA, and KCLUS. The main content area displays the title 'SAS Visual Statistics' and a brief description: 'SAS Visual Statistics provides access to SAS/STAT procedures and SAS/GRAPH procedures.'

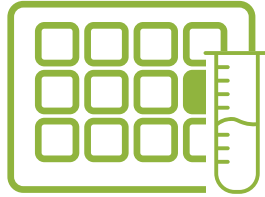
- [SAS Viya: Machine Learning](#)

The screenshot shows the SAS Viya Platform Programming Documentation interface. The top navigation bar includes 'SAS Viya Platform Programming Documentation | 2023.03'. The left sidebar contains a navigation menu with categories like 'Welcome to SAS Programming Documentation', 'What's New', 'Learning SAS Viya Platform Programming', 'Syntax Quick Links', and 'SAS Procedures by Name and Product'. The 'SAS Procedures by Name and Product' section is expanded, showing a list of procedures: ASTORE, BNET, BOOLRULE, FACTMAC, FASTKNN, FISM, FOREST, GMM, and GRADBOOST. The main content area displays the title 'SAS Viya: Machine Learning' and a brief description: 'SAS Viya: Machine Learning provides access to SAS Viya Machine Learning procedures and SAS Viya Machine Learning procedures.'

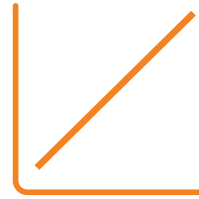
Data Preparation



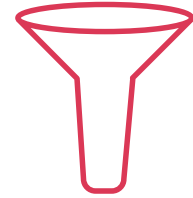
binning



sample



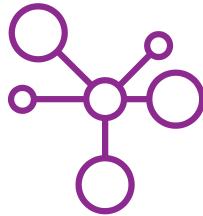
transformation



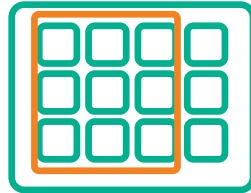
filter



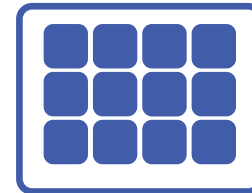
re-scaling



variable
clustering



custom
subset
selection

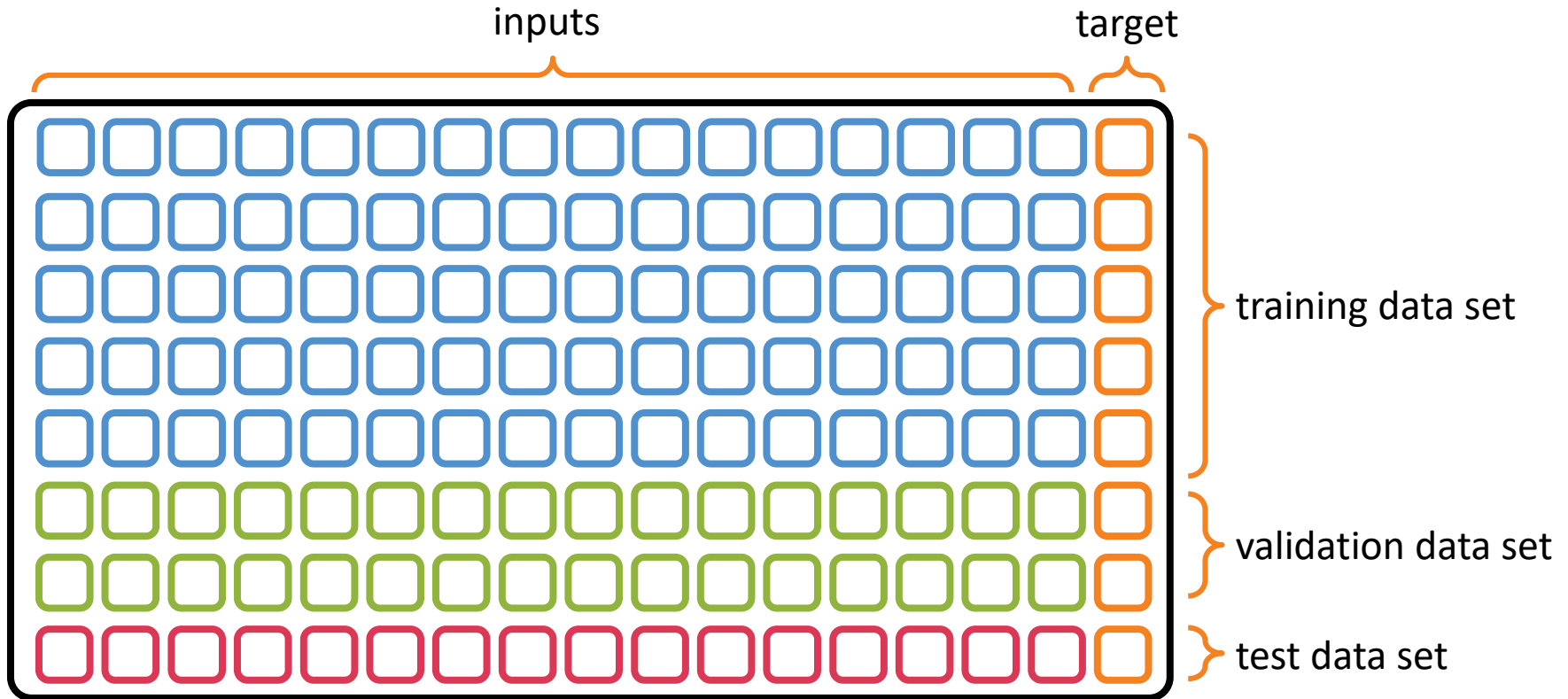


imputation

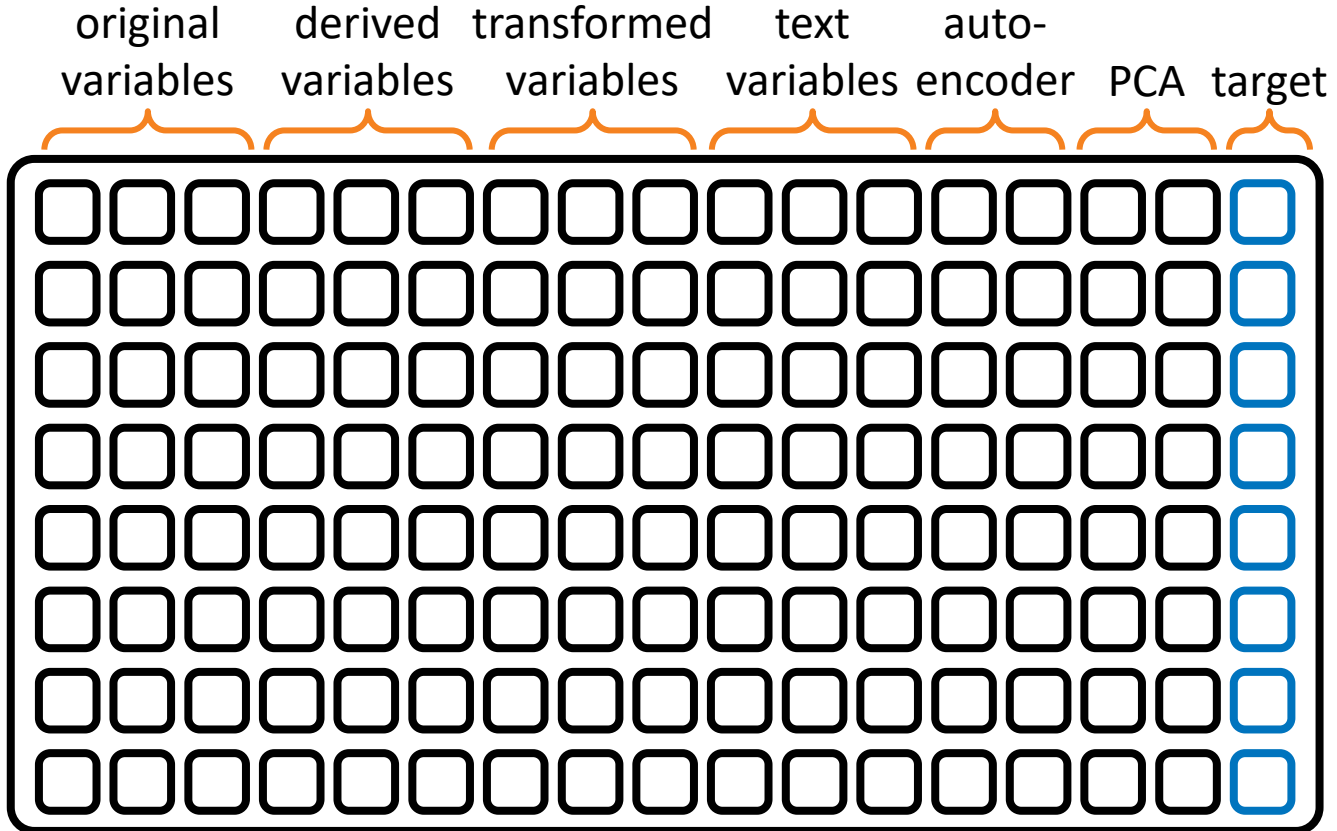


feature
extraction

Partitioning

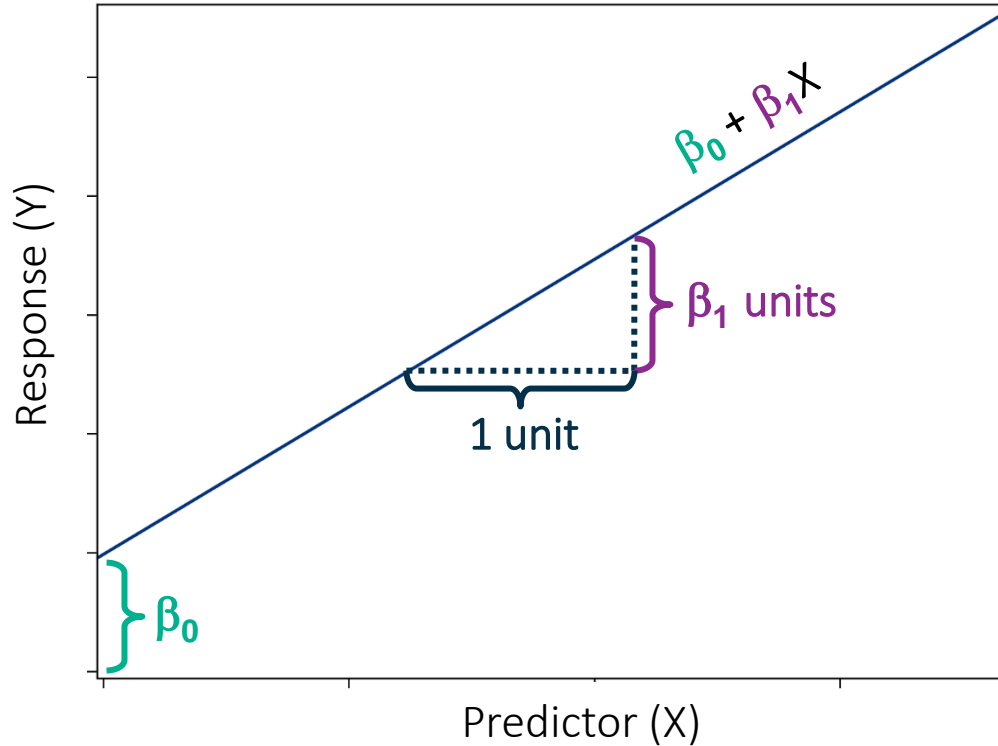


Feature Extraction

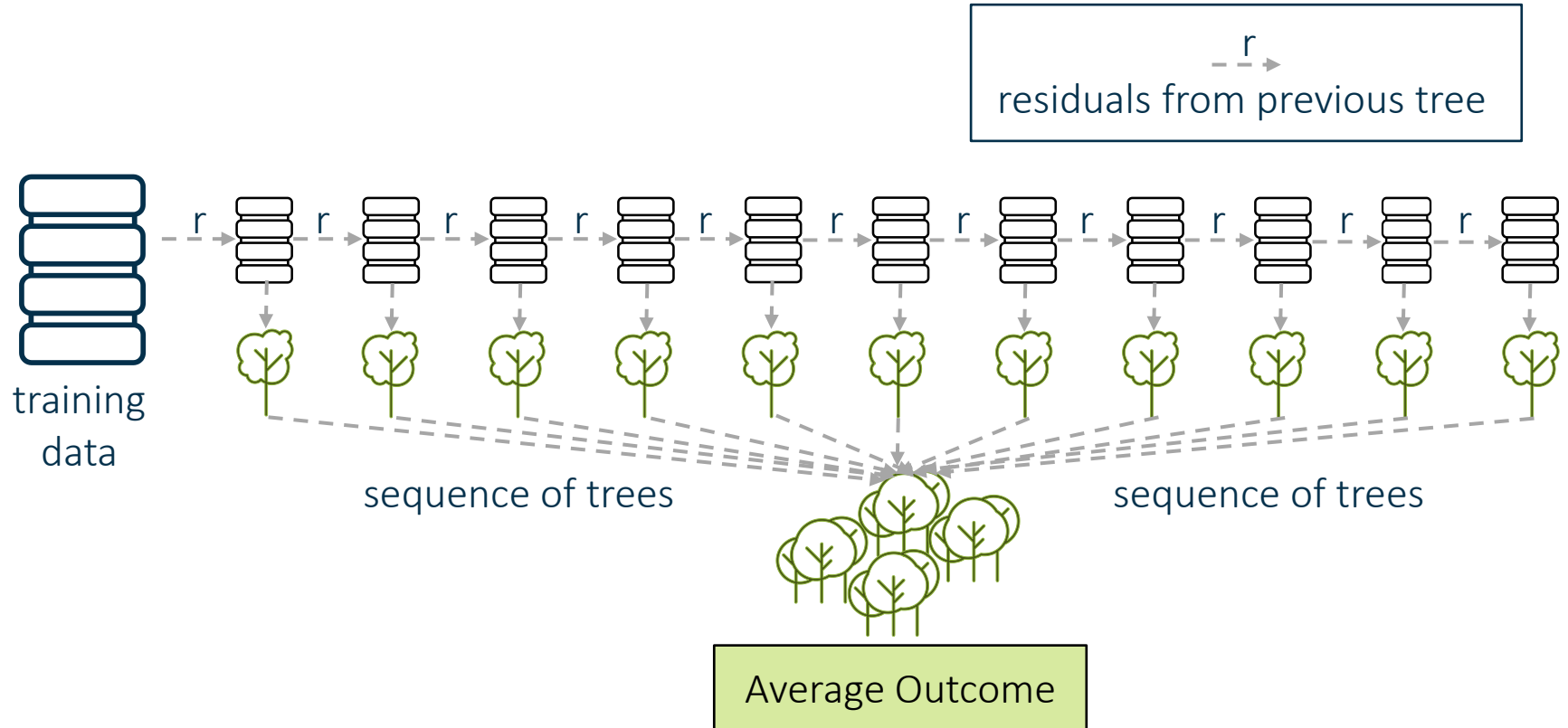


Ordinary Least Squares (OLS) Regression Model

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$



Ensemble of Trees: Gradient Boosting



Lift Value

- Multiplicative Factor how your predictive model is better than a random selection.

