

머글들을 위한 Machine Learning 3

김영욱

부장 / PS / Microsoft

youngwook@outlook.com

Blog: Youngwook.com

문제 정의

어떤 모델을 만들 것인가?



미국 45개 월마트 부서별 주간 판매액 예측 모델

kaggle™



Competitions

Documentation

InClass

General

InClass

Sort by Grouped

All Categories

Search competitions



19 Active Competitions



Severstal: Steel Defect Detection

Can you detect and classify defects in steel?

Featured - Code Competition - a month to go - manufacturing, image data

\$120,000

1,665 teams



Lyft 3D Object Detection for Autonomous Vehicles

Can you advance the state of the art in 3D object detection?

Featured - 2 months to go - image data, object detection

\$25,000

123 teams



RSNA Intracranial Hemorrhage Detection

Identify acute intracranial hemorrhage and its subtypes

Featured - 2 months to go - health foundations and medical research, image data

\$25,000

399 teams



The 3rd YouTube-8M Video Understanding Challenge

Temporal localization of topics within video

Research - 16 days to go - video data, object detection

\$25,000

263 teams

<https://www.kaggle.com/c/walmart-recruiting-store-sales-forecasting/data>



Walmart Recruiting - Store Sales Forecasting

Use historical markdown data to predict store sales

690 teams · 5 years ago

[Overview](#) [Data](#) [Notebooks](#) [Discussion](#) [Leaderboard](#) [Rules](#)

[Join Competition](#)

Data Description

You are provided with historical sales data for 45 Walmart stores located in different regions. Each store contains a number of departments, and you are tasked with predicting the department-wide sales for each store.

In addition, Walmart runs several promotional markdown events throughout the year. These markdowns precede prominent holidays, the four largest of which are the Super Bowl, Labor Day, Thanksgiving, and Christmas. The weeks including these holidays are weighted five times higher in the evaluation than non-holiday weeks. Part of the challenge presented by this competition is modeling the effects of markdowns on these holiday weeks in the absence of complete/ideal historical data.

stores.csv

This file contains anonymized information about the 45 stores, indicating the type and size of store.

train.csv

데이터 전처리 과정

데이터 지원 형식 확인 및 변환

데이터 업로드

데이터 Merge

데이터 전처리

데이터 지원 형식

- 헤더가 있거나 (.csv) 없는 (.nh.csv) 쉼표로 구분된 값(csv)
- 헤더가 있거나 (.tsv) 없는 (.nh.tsv) 탭으로 구분된 값 (tsv)
- 일반 텍스트(txt)
- Excel file, Azure Table, Hive Table
- SQL Server
- SVMLight 데이터(.svmlight)
- 특성 관계 파일 형식(ARFF) 데이터(. Arff)
- Zip 파일(.zip)
- R 개체 또는 작업 영역 파일(.RData)

데이터 지원 유형

- 문자열
- 정수
- Double
- Boolean
- Datetime
- timespan

데이터 지원 용량

- 1.98 GB/file
- Total 10GB

45개 Walmart 부서별 주간 판매량

2010-02-05 ~ 2012-11-01

12열 X 8,190행

○ feature.csv

● 지점

● 날짜

● 온도

● 연료비

● 프로모션 * 5

● 소비자 물가 지수

● 실업률

● 휴일 여부

5열 X 421,570행

○ train.csv

● 지점

● 부서

● 날짜

● 판매량

● 휴일 여부

3열 X 45행

○ stores.csv

● 지점

● 유형

● 규모

실습1:

Walmart data preprocessing

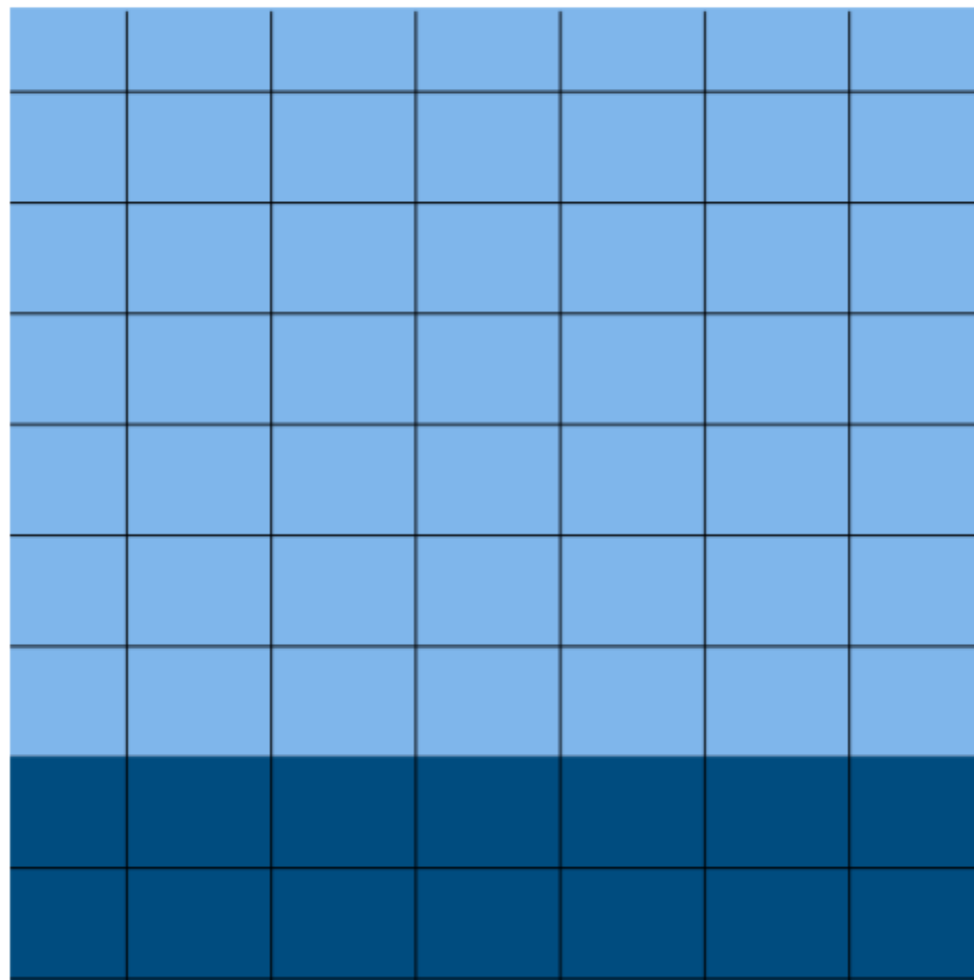
$$\mathcal{X} = \text{Feature}$$
$$\mathcal{Y} = \text{Label}$$
[illegible]

\mathcal{X} = *Feature*

\mathcal{Y} = *Label*

X_{train}

X_{test}



Y_{train}

Y_{test}



Random Split

\mathcal{X} = *Feature*

\mathcal{Y} = *Label*

True
True
True
True
True
False
False
False
False

X_{train}

X_{test}

True

True

True

True

True

False

False

False

False

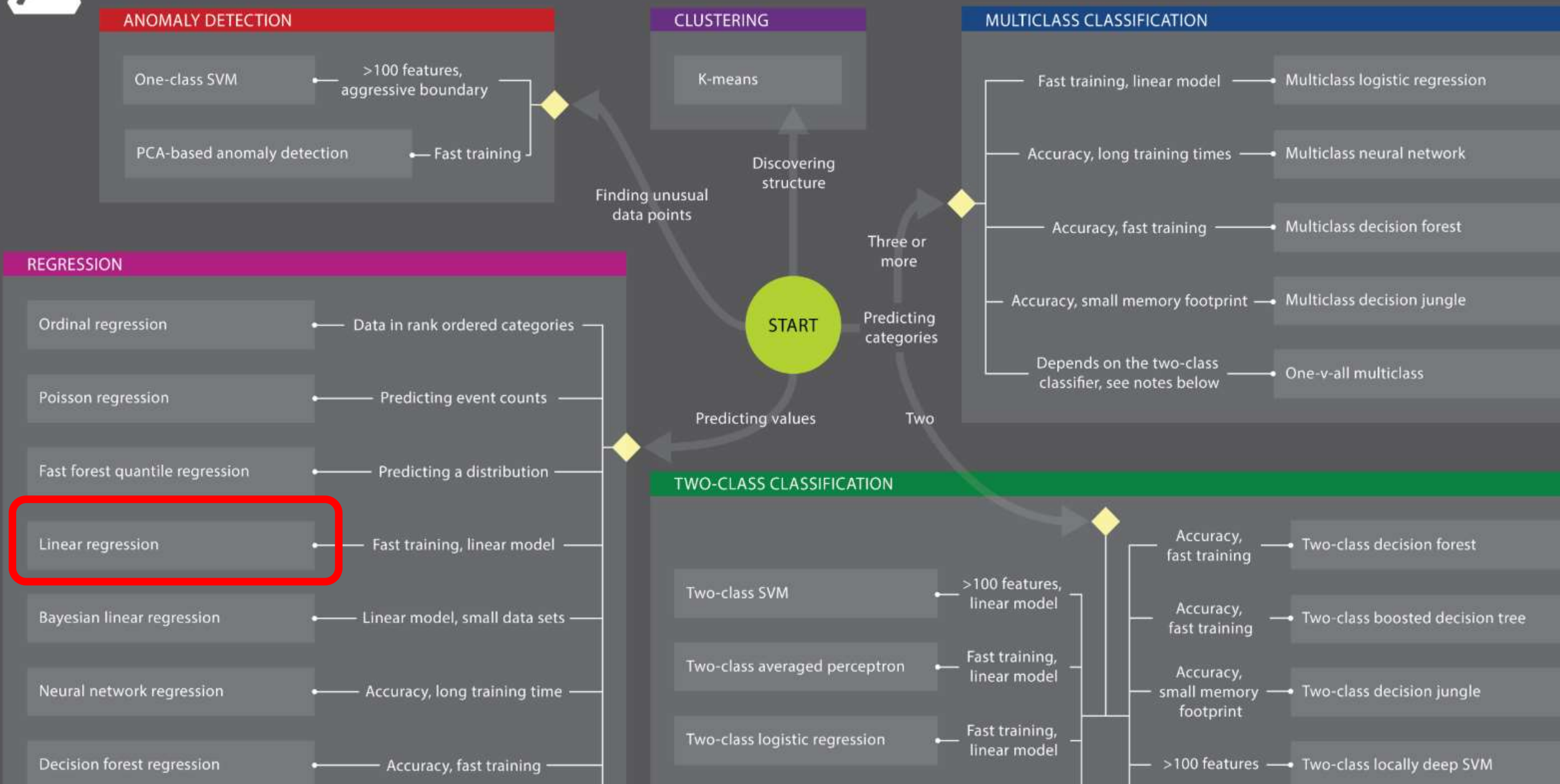
Y_{train}

Y_{test}



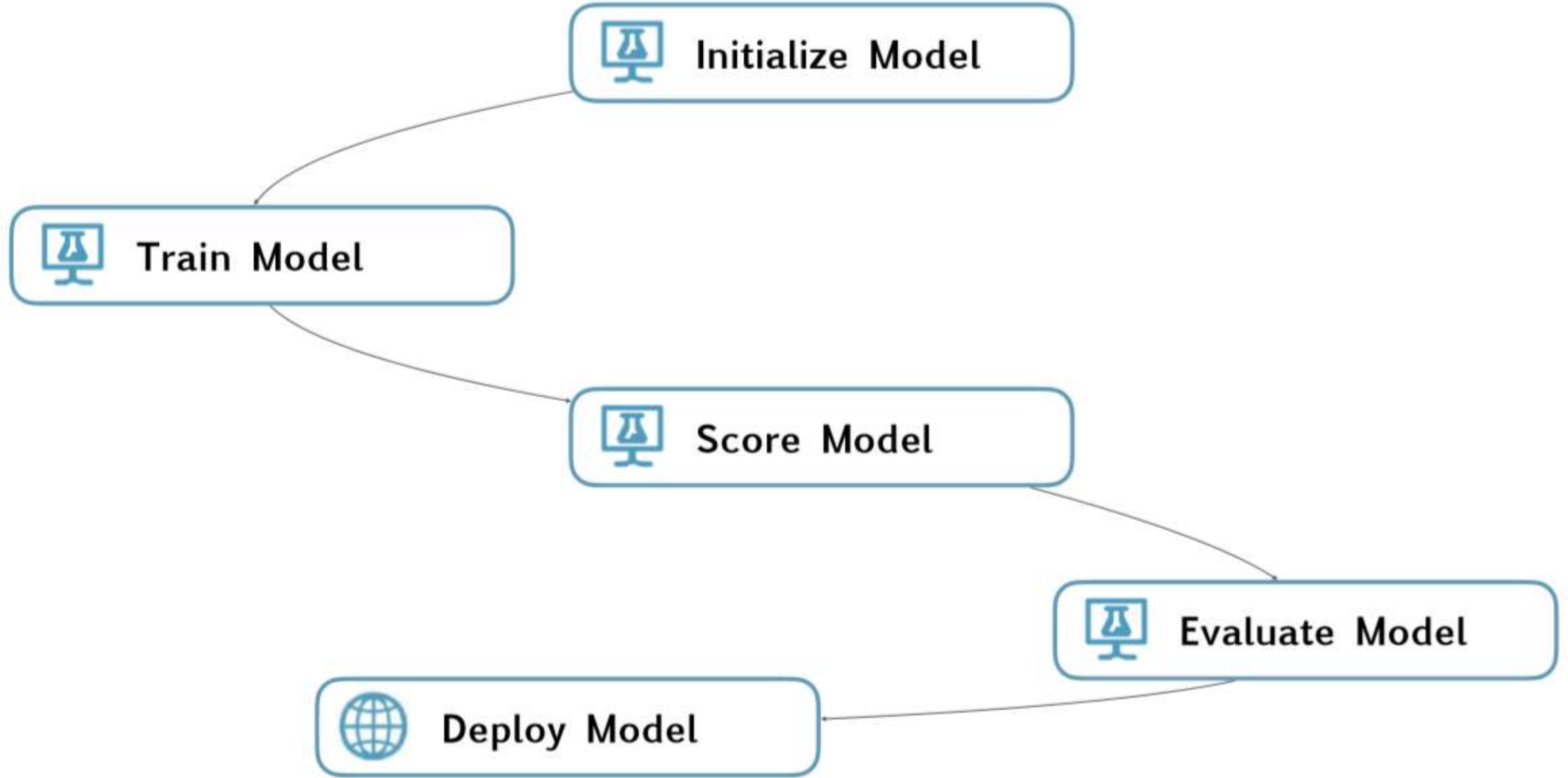
Microsoft Azure Machine Learning: Algorithm Cheat Sheet

This cheat sheet helps you choose the best Azure Machine Learning algorithm for your predictive analytics solution. Your decision is driven by both the nature of your data and the question you're trying to answer.



실습1:

Walmart Machine Learning



Microsoft Professional Program

감사합니다
Thank you~!