

Final Project

- Data collection / exploration / visualization
- Data preprocessing & Feature engineering
- Train & test traditional ML algorithms
- Train & test deep-learning models
- Compare various models & deliver the result

수업 관련 공지사항

* 데이터 선정 / 전처리 프로세스 / 모델 & 지표 선택 모두 자유입니다. (배운 내용의 복습에 Focus!)

* Part 1~6 에서 배운 지식들을 최대한 모두 활용하는데 초점을 맞춰주세요. (크롤링 필수 X)

- 12/3(목)~12/4(금) : 문제 정의 / 데이터 선정 / 데이터 탐색 & 시각화
- 12/7(월)~12/9(수) : 데이터 전처리 / ML & DL 모델 적용 / 모델 튜닝 & 최종 모델 선택
- 12/10(목) : 발표 준비 / 최종 발표

- 12/10 (목) 3~4교시 : 팀별 발표 및 질의응답 (15~20분 내외/팀)

: 목요일 2교시 종료 전까지 발표 자료 & Jupyter notebook 제출 : 슬랙 DM or repositivator@gmail.com

* 발표 시작 시간은 일정에 따라 변동될 수 있습니다 & 도움이 필요할 경우 슬랙 채널에서 호출

Various data collection – etc (Datasets / Data repository)

Awesome Public Datasets @ <https://github.com/awesomedata/awesome-public-datasets>

Google AI Datasets @ <https://ai.google/tools/datasets>

Google Dataset Search @ <https://toolbox.google.com/datasetsearch>

SKT BigData Hub @ <https://www.bigdatahub.co.kr>

Kaggle competition datasets @ <https://www.kaggle.com/datasets>

(ex. Google Play Store Apps data @ <http://j.mp/2PDhbKR>)

<http://www.aihub.or.kr> – AI 오픈이노베이션 허브 (한국어 음성 & 대화, 한국인 안면, 법률/특허/헬스케어/관광/농업/이미지 데이터)

<https://golmok.seoul.go.kr> – 서울시 우리마을가게 상권분석 서비스

<http://data.seoul.go.kr> – 서울 열린 데이터 광장

<https://www.dataquest.io/blog/free-datasets-for-projects> – 19 Places to Find Free Data Sets for Data Science Projects

<http://dataportals.org> – A Comprehensive List of Open Data Portals from Around the World

<https://www.kdnuggets.com/datasets/index.html> – Datasets for Data Mining/Science

<http://figshare.com> – Help academic institutions store, share and manage their research

<https://opendatainception.io> – 2600+ Open Data Portals around the World

<https://search.datacite.org> – Locate, identify, and cite research data

<http://aws.amazon.com/datasets> – AWS dataset

<http://quandl.com> – Financial Data

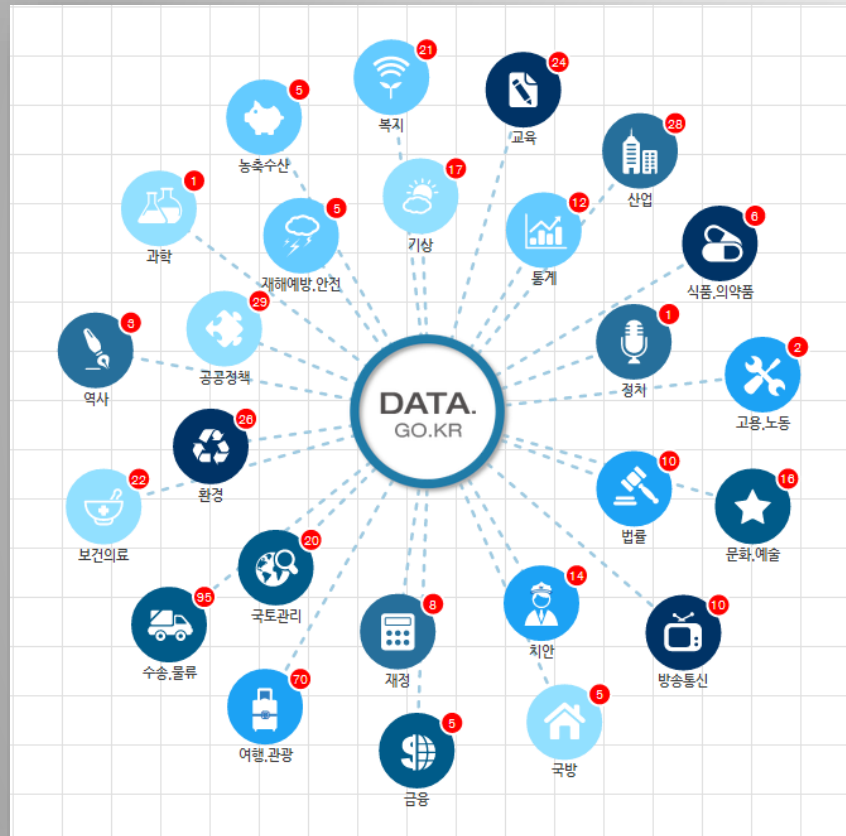
* 각종 데이터분석 관련 공모전/대회/프로젝트사례 모음 @ <http://j.mp/2MPDfON>

* 해외 기업의 인공지능 데이터 개방과 활용 현황 (구글 사례를 중심으로) @ <http://j.mp/2paKt7j>

* 딥러닝 학습을 위한 국내외 데이터셋 현황 (이미지/동영상/음성) @ <http://j.mp/2BSShNy> / <http://j.mp/2roFj8i> / <http://j.mp/2VbHXeG>

본 교안은 광주 ICT 교육 과정을 위해 제작되었으며, 교육 외 배포/게시/공개를 금합니다.

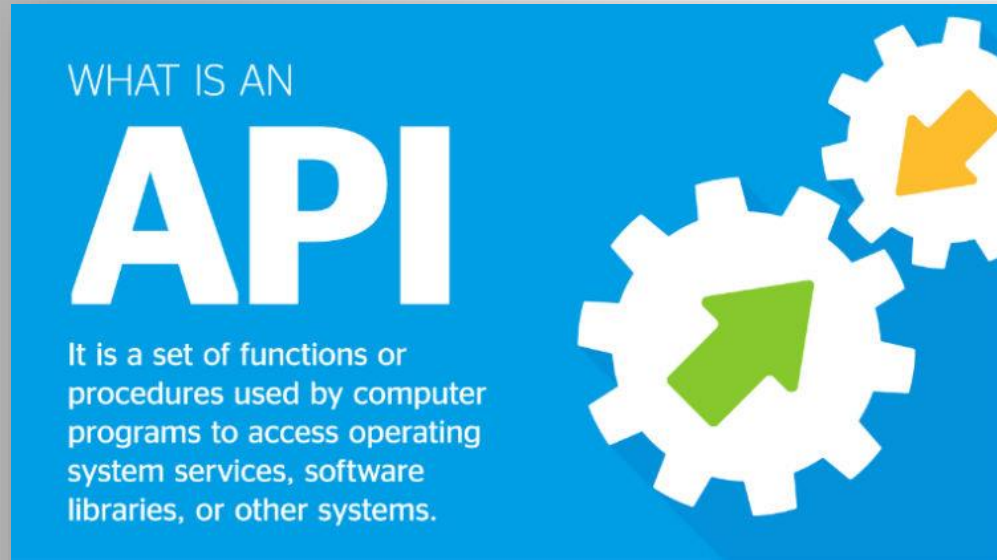
Various data collection – Public data & Open data (APIs & files)



- 공공 데이터 포털 : <https://www.data.go.kr>
- 국가 통계 포털 : <http://kosis.kr>
- MDIS (MicroData Integrated Service) : <https://mdis.kostat.go.kr>

1. Available resources for data collection

Various data collection – Unowned data



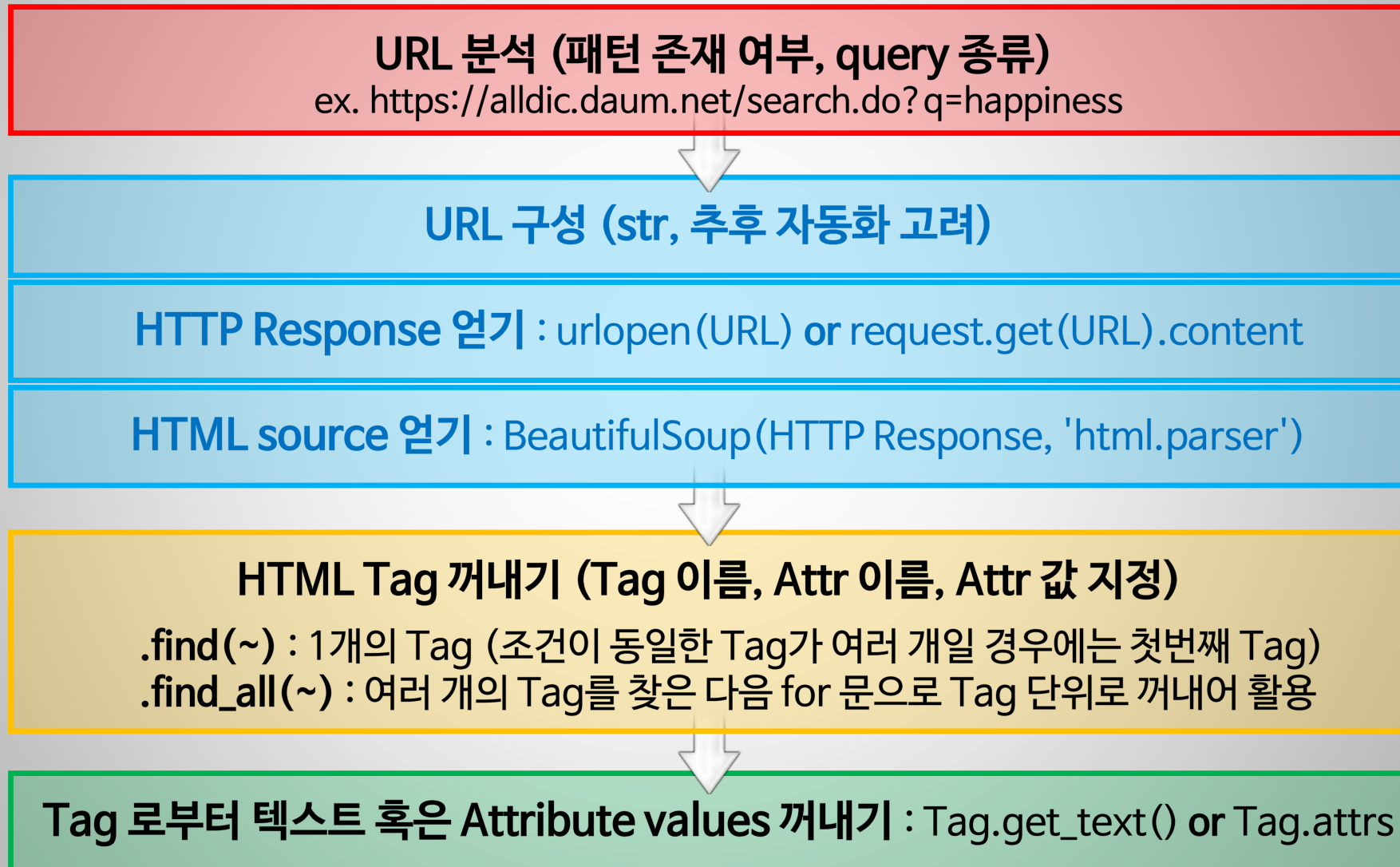
Use APIs & Web scraper

- APIs (Twitter, Facebook, Instagram, etc)
- Bots (Web crawler, Web scraper)

* 네이버 크롤링 라이브러리 Kocrawl (날씨/미세먼지/지도/맛집/맞춤법) @ <https://j.mp/2CbdRA8>
* 여기어때, 야놀자 DB 무단수집 위법 판결 @ <https://j.mp/3fgxi9Q> + 크롤링과 저작권 침해 고소 진행 일대기 @ <https://j.mp/3k5vbbl>
* Web Scraping Tool & Web Data Extractor : ScrapeStorm (\$49/month) @ <http://j.mp/2Y4porj> + Octoparse @ <https://j.mp/3iYQTgX>

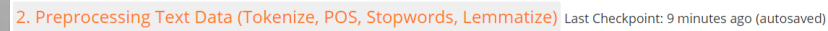
본 교안은 광주 ICT 교육 과정을 위해 제작되었으며, 교육 외 배포/게시/공개를 금합니다.

The process of web scraping (detailed)



1. Available resources for data collection

The process of data analysis for text data



File Edit View Insert Cell Kernel Help

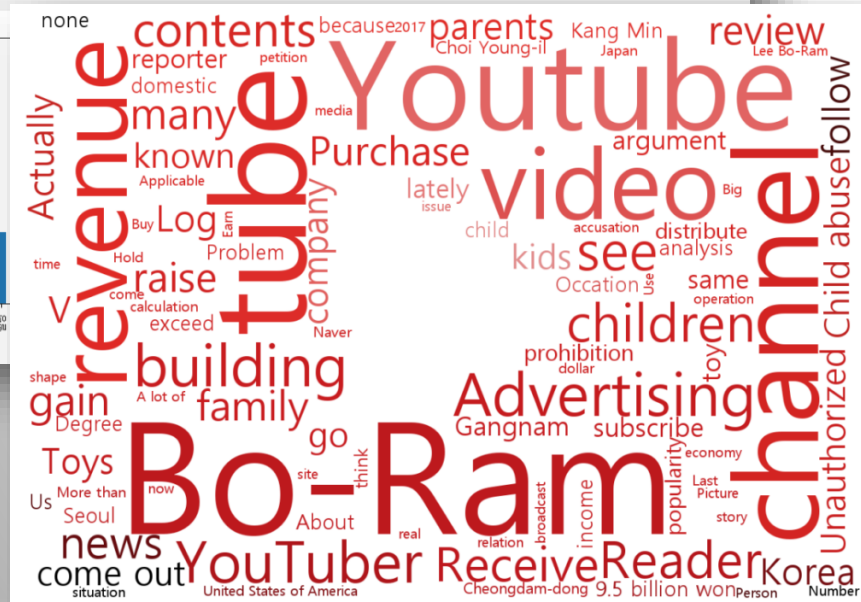
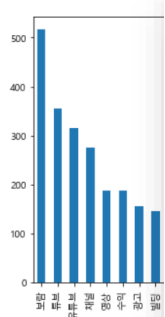
Text data Preprocessing

- nltk library(Natural Language Toolkit)를 이용하여 Text Processing을 위한 전처리를 실습한다.

1. 영어 문장 토큰화하기

```
In [ ]: 1 # Test processing을 위해 nltk package 를 import!  
2 import nltk
```

```
In [ ]: 1 #
2 #
3 #
4 #
5 #
6 plt.show()
```



텍스트 데이터를 str 자료형으로 준비

Tokenize (형태소 분석)

POS Tagging (Part-of-speech, 품사 표시)

Stopwords 제거 (불용어 제거)

단어 갯수 카운팅 & 단어 사전 생성

단어 사전 기반 데이터 시각화

(+ 머신러닝/딥러닝 모델 적용)

2. Possible pathways for data preprocessing

- + Check & adjust data for handling **Missing data & Outlier**
- + Select important columns (or just use all columns & improve your model later)
- + Change characters to numbers (Binary num, Class num, One-hot vector, etc.)
- + (If applicable & useful) **Select features** with Tree-based models
- + (If applicable & useful) **Modify the scale of features** with StandardScaler / MinMaxScaler
- + (If applicable & useful) **Reduce dimension** with PCA
- + (If applicable & useful) Try **other traditional ML Models** for enhancing the result (except DL/NN)

서울시 범죄현황 통계자료 분석 및 시각화

2. 서울시 범죄현황 통계자료 분석 및 시각화 Last Checkpoint: an hour ago (unsaved changes)

File Edit View Insert Cell Kernel Help Trusted

서울시 범죄현황 통계자료 분석 및 시각화

In [1]:
1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4
5 import matplotlib.pyplot as plt
6 from matplotlib import font_manager, rc

1. 데이터 입력 및 데이터 전처리

In [2]:
1 df = pd.read_excel('관서별 5대범죄 발생 및 검거.xlsx', encoding='utf-8')
2 df.head()

	관서명	소계(발생)	소계(검거)	살인(발생)	살인(검거)	강도(발생)	강도(검거)	강간(발생)	강간(검거)	절도(발생)	절도(검거)	폭력
0	계	126401	82688	163	156	276	257	5449	5069	55307	21842	652
1	충무서	2860	1716	2	2	3	2	185	65	1395	477	135
2	종로서	2472	1589	3	3	6	5	115	98	1070	413	127
3	남대문서	2094	1226	1	0	6	4	65	46	1153	382	869
4	서대문서	4029	2579	2	2	5	4	154	124	1812	738	205

Scikit-learn practices & Appendix

(Appendix 0) Missing data visualization (with missingno).zip

(Appendix 1) Tuning HyperParams with Hyperopt (+ LightGBM).zip

(Appendix 2) Auto-ScikitLearn (with Breast cancer data).zip

(Appendix 3) Auto-feature-engineering with FeatureTools.zip

(Appendix 4) PCA for BreastCancer & Cifar10 (딥러닝 학습 후 추가 학습).zip

(Appendix 5) Colab 기준 LightGBM 설치 방법 (2020.03).zip

1. (Cheat Sheet) Scikit_Learn.zip

2. Hands On MachineLearning with ScikitLearn and TensorFlow (1판 & 2판).zip

3. Model saving & loading (Scikit-learn).zip

4. Amazon SageMaker Hands-on.docx

5. AI Career Pathways with Roles, Tasks, Skills (Workera Report).pdf

파이썬을 활용한 기초 통계분석

2. 빈도 분석 & 기술통계량 분석

File Edit View Insert Cell Kernel Widgets Help

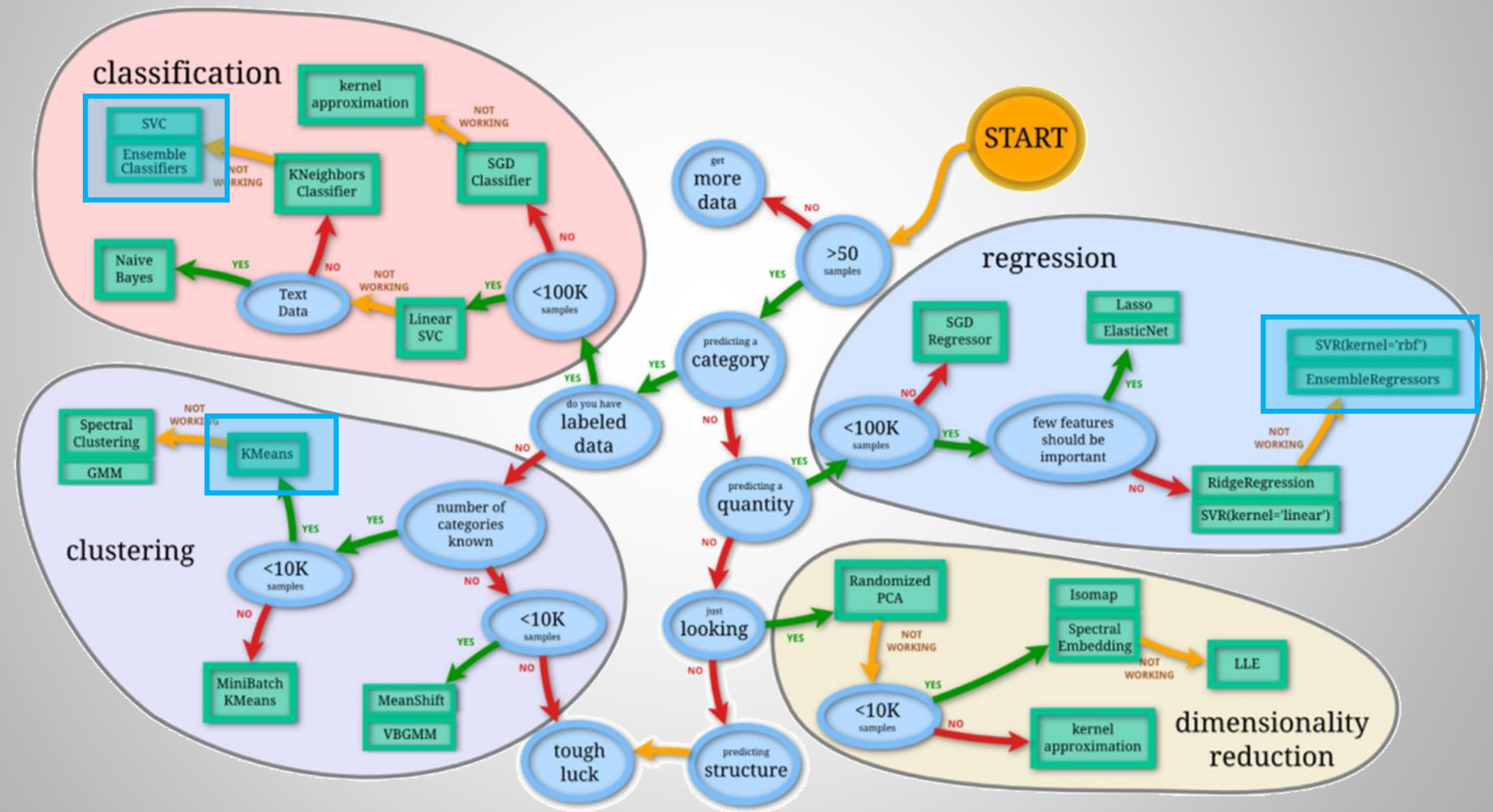
3. Outlier 의 탐지 및 제거와 전후 분포 비교

In [210]:
1 df.boxplot(column='amount') # 위 아래의 작은 점(검정색 선이 상/하한선, 그 밖의 경우는

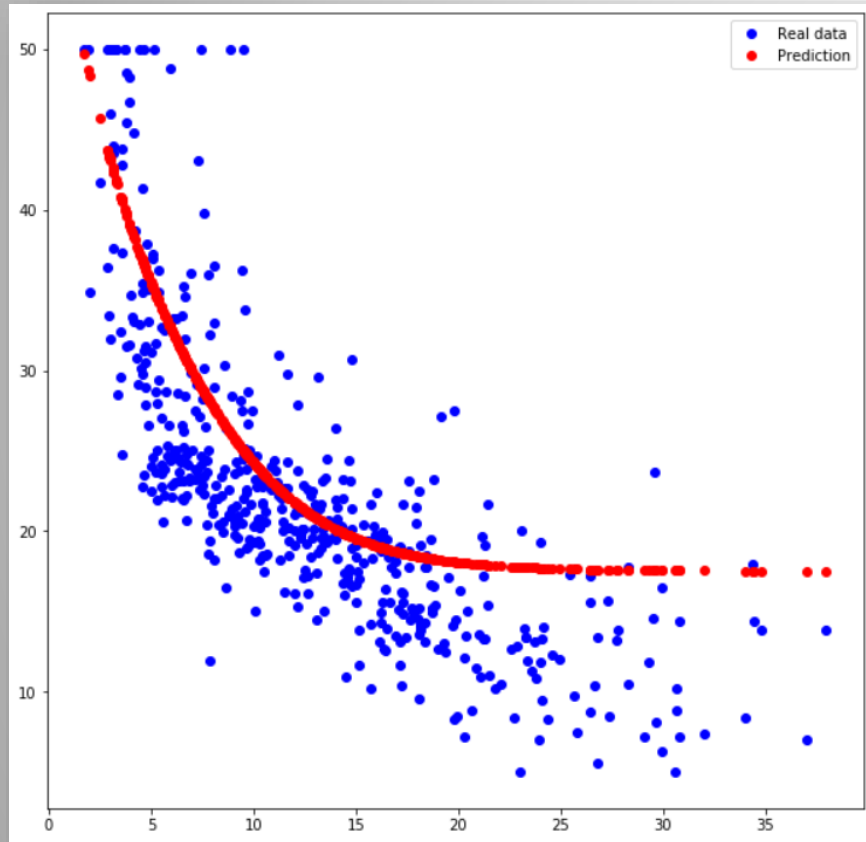
<matplotlib.axes._subplots.AxesSubplot at 0x1b1170839b0>

풀어내려는 문제의 종류와 데이터의 타입(형태, 수)에 따른 ML 알고리즘 선택 가이드

http://scikit-learn.org/stable/tutorial/machine_learning_map/ (각 알고리즘 별 예시 코드 有)



Neural-network modeling with TensorFlow & Keras



- 0-1. (UseThis) Classification (Titanic dataset).ipynb
- 0-2. (UseThis) Classification with Keras (Titanic dataset).ipynb
- 0-3. (UseThis) Regression (Boston house price dataset).ipynb
- 0-4. (UseThis) Regression with Keras (Boston house price dataset).ipynb

Try other improvements,

- Other **activation functions** (tanh, relu)
- Other **optimizers** (Adam, Adagrad, RMSProp)
- Other **learning rates** (0.01, 0.0001)
- More **learning steps** (75000, 100000)
- More **layers & nodes** (64, 128, 256)
- + **Model stacking**
- + **AutoML** (Keras-tuner, Google AutoML Tables, FeatureTools 등)
- + **Bayesian Hyperparams Optimization**

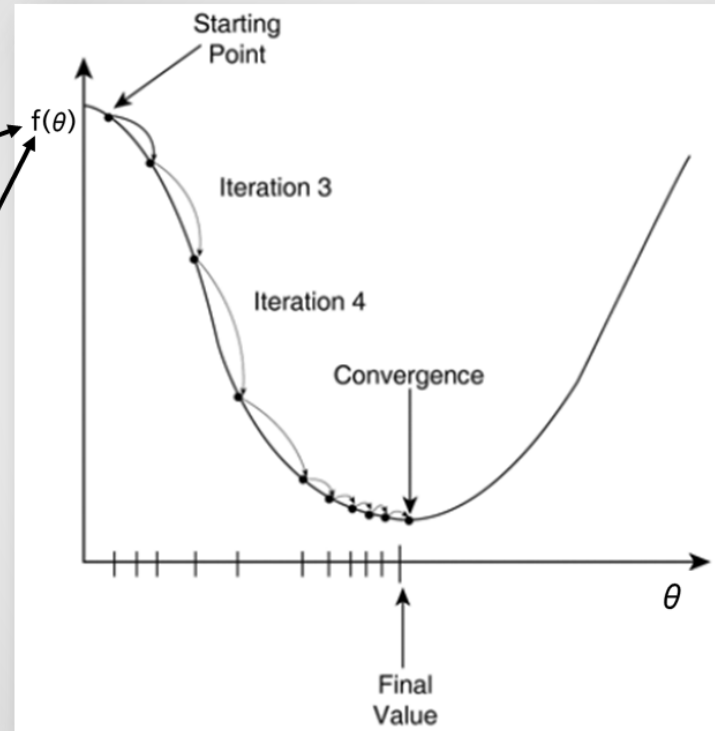
4. Test the model with appropriate metrics

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

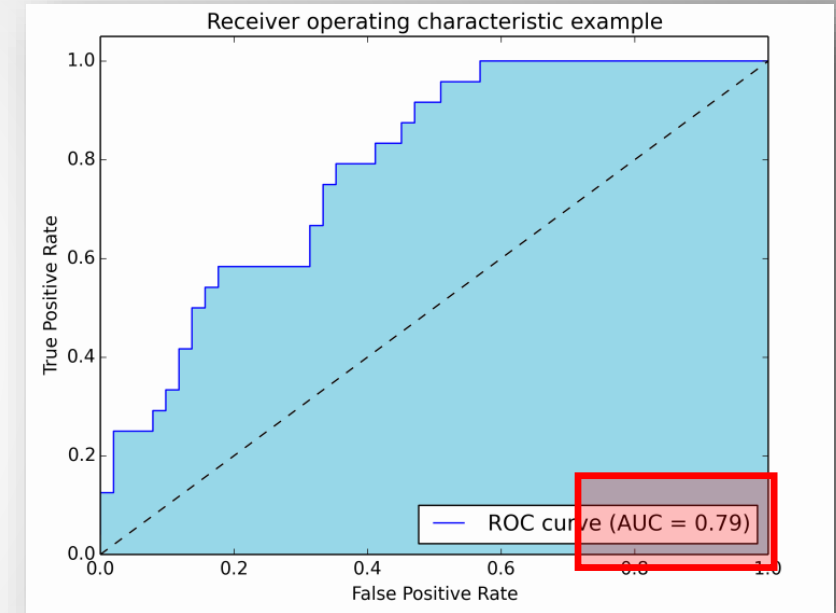
Mean squared error
for regression

$$J(\theta) = - \sum_i y^{(i)} \log(h_{\theta}(x^{(i)}))$$

Cross-entropy
for classification



AUC = Area Under the ROC Curve



- measures the **quality** of classifier.
- AUC = 0.5 : random classifier.
- AUC = 1 : **perfect** classifier.

* 발표 시 포함할 사항 :

1. 프로젝트 소개 (어떤 분석을 하였는가)
2. 데이터 소개/탐색/시각화 (출처, 형식, 분포 등)
3. 데이터 전처리 과정 (적용한 전처리 방법 & 이유)
4. 적용한 분석 기법 및 모델 소개
5. 분석 및 모델링 결과 (각종 지표 수치 제시)
(+ 가능 시 추가로 분석하면 좋을 과제 제시)

* 발표 시 제출할 사항 :

전체 코드 with 주석
(.ipynb, 단일 혹은 복수)

발표 자료 필수
(PPT or PDF)

수업 관련 공지사항

1팀 : 한대영, 고성현, 문수지, 서경호

2팀 : 강예은, 안현진, 조우람, 최갑주

3팀 : 박광렬, 손부언, 이승재, 이하은

4팀 : 김관주, 강기훈, 조혜빈, 정경훈

광주 ICT 이노베이션스퀘어
인공지능 융복합교육

End of Document