

프로젝트 기반 데이터 과학자 양성과정(Data Science) Machine Learning 및 분석실습

8주차
지도 학습

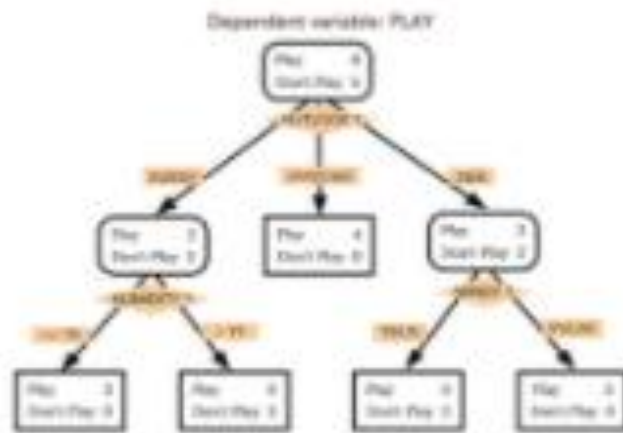
Gradient Boosting

강사 : 최영진

1. Gradient Boosting

❖ Random Forest

- 의사결정 트리의 오버피팅 한계를 극복하기 위한 전략으로 **랜덤 포레스트(Random Forest)** 등장
- 데이터에 의사결정나무 여러 개를 동시에 적용해서 학습성능을 높이는 앙상블 기법
- 동일한 데이터로부터 복원추출을 통해 30개 이상의 데이터 셋을 만들어 각각에 의사결정나무를 적용한 뒤 학습 결과를 취합하는 방식



Tree

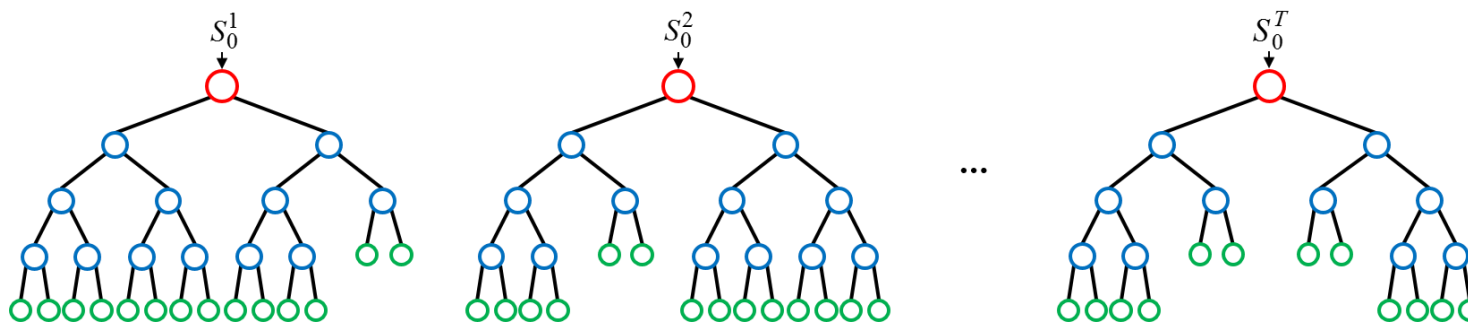


Random Forest

1. Gradient Boosting

❖ Random Forest

- 배깅(bagging): bootstrap aggregating의 약자로, 부트스트랩(bootstrap)을 통해 조금씩 다른 훈련 데이터에 대해 훈련된 기초 분류기(base learner)들을 결합(aggregating)시키는 방법
- bootstrap sampling(복원추출)을 사용하며 decision tree 생성으로 algorithm으로 진행
- 트리들의 편향은 그대로 유지하면서, 분산은 감소시키기 때문에 포레스트의 성능을 향상



부트스트랩 방법을 통해 T 개의 훈련 데이터셋을 생성
 T 개의 기초 분류기(트리)들을 훈련
기초 분류기(트리)들을 하나의 분류기(랜덤 포레스트)로 결합(평균 또는 과반수투표 방식 이용).

1. Gradient Boosting

❖ Random Forest

▪ 장점

- 다양성을 극대화 하여 예측력이 상당히 우수한 편
- 다수의 트리의 예측 결과를 종합하여 의사결정을 진행하기 때문에 안정성도 상당히 높음
- 랜덤화(randomization)는 포레스트가 노이즈가 포함된 데이터에 대해서도 강인

▪ 단점

- 다수의 트리를 이용한 의사결정 기법을 이용하기 때문에 기존의 트리가 갖는 장점 중 하나인 설명력을 잃음

1. Gradient Boosting

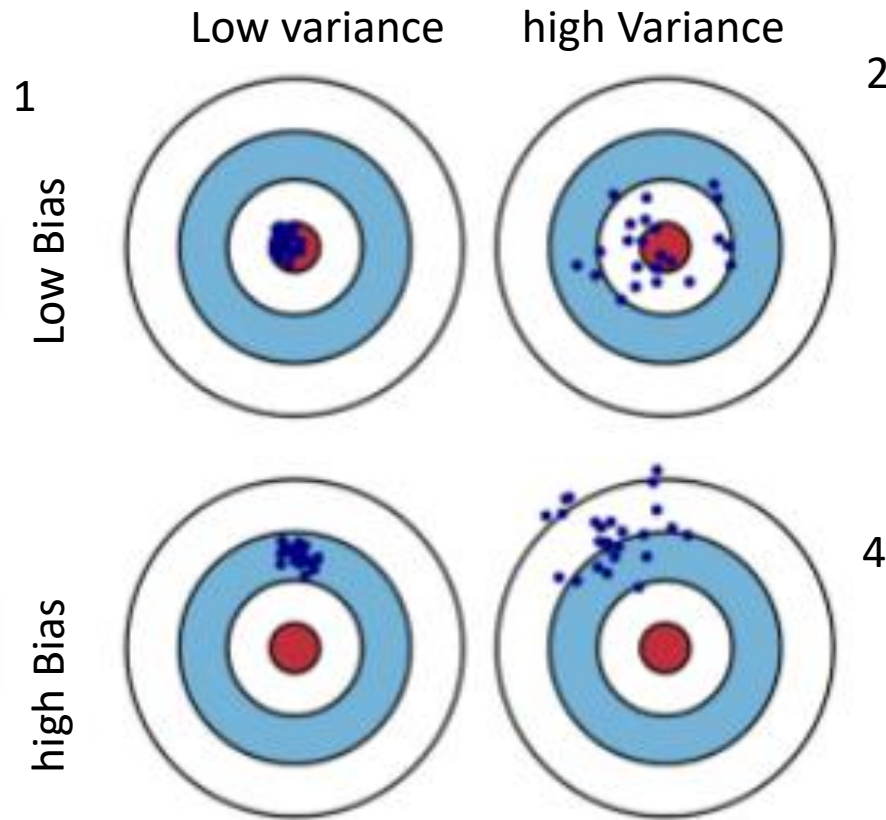
❖ Ensembles(앙상블)

- Error 최소화
 - 다양한 모델의 결과를 종합하여 데이터의 오류를 줄여줌
- Overfitting 감소
 - 각 모델별로 bias가 존재
 - Bias를 종합하여 결과를 생성하고, overfitting을 줄임
- **bias와 variance는 모델의 loss 또는 error를 의미**

1. Gradient Boosting

❖ Ensembles(앙상블)

- bias와 variance는 모델의 loss 또는 error를 의미



1번 과적

예측값들이 대체로 정답 근방 → 편향이 낮음
예측값들끼리 서로 몰릴 경우 → 분산이 낮음

2번 과적

예측값들이 대체로 정답 근방 → 편향이 낮음
예측값들끼리 서로 흩어져 있는 경우 → 분산이 높음

3번 과적

예측값들이 대체로 정답으로부터 멀어져 있는 경우 → 편향이 높음
예측값들끼리 서로 몰려 있는 경우 → 분산이 낮음

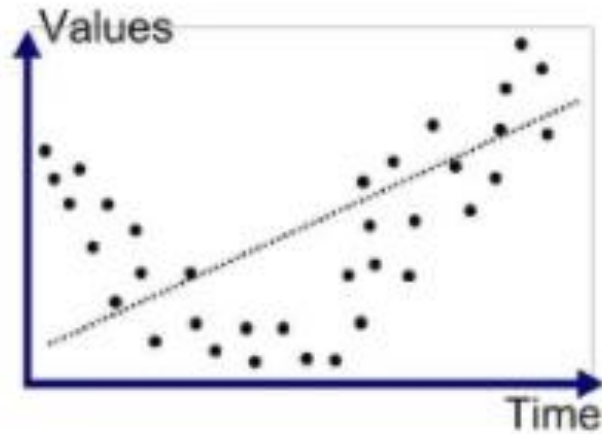
4번 과적

예측값들끼리 대체로 정답으로부터 멀어져 있는 경우 → 편향이 높음
예측값들끼리 서로 흩어져 있는 경우 → 분산이 높음

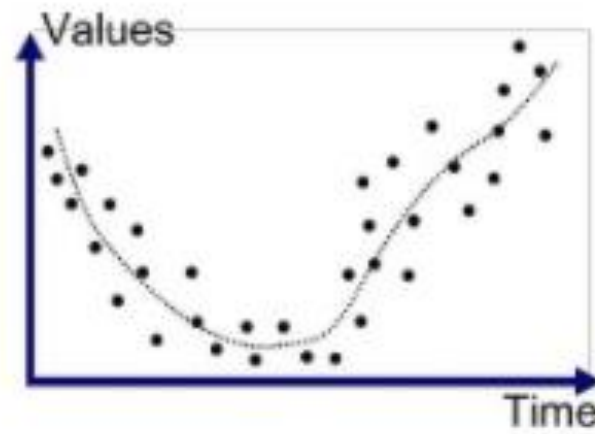
1. Gradient Boosting

❖ Ensembles(앙상블)

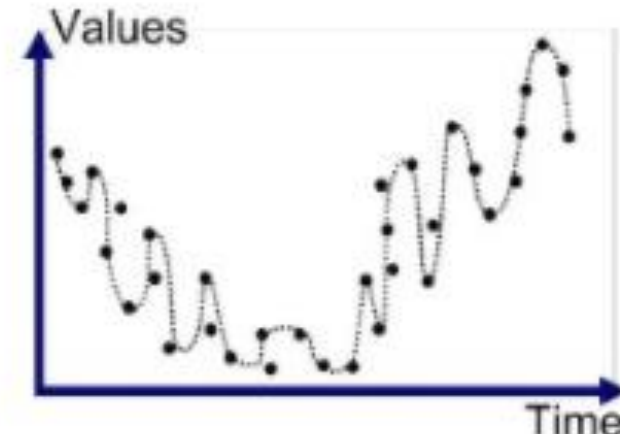
- 예측값들과 정답이 대체로 멀리 떨어져 있으면 결과의 편향(bias)이 높음
예측값들이 자기들끼리 대체로 멀리 흩어져있으면 결과의 분산(variance)이 높음



high bias
low variance



medium bias
medium variance

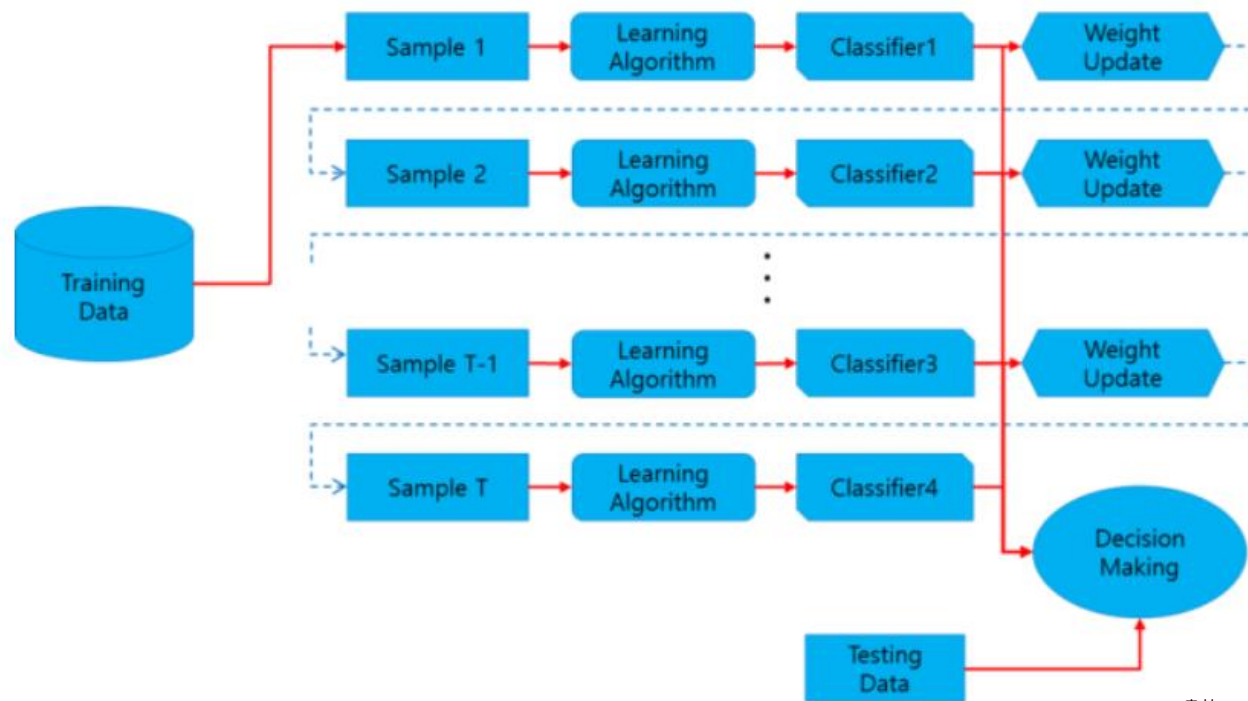


low bias
high variance

1. Gradient Boosting

❖ Boosting이란?

- Bagging과 유사하게 초기 샘플 데이터를 조작하여 다수의 분류기를 생성하는 기법 중 하나지만 가장 큰 차이는 순차적(Sequential)방법
- 이전 분류기의 학습 결과를 토대로 다음 분류기의 학습 데이터의 샘플가중치를 조정해 학습을 진행하는 방법

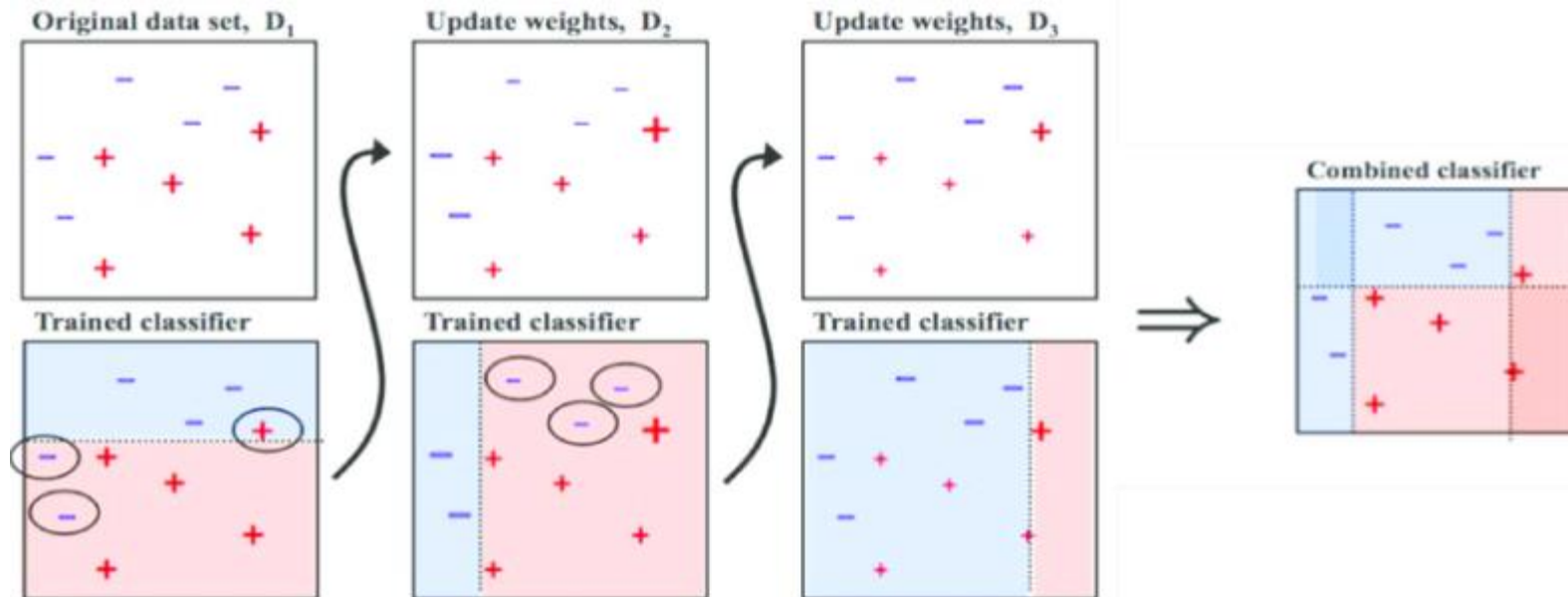


[출처] [EM] 부스팅(Boosting) 이론

1. Gradient Boosting

❖ Boosting이란?

- Bagging과 유사하게 초기 샘플 데이터를 조작하여 다수의 분류기를 생성하는 기법 중 하나지만 가장 큰 차이는 순차적(Sequential)방법
- 이전 분류기의 학습 결과를 토대로 다음 분류기의 학습 데이터의 샘플가중치를 조정해 학습을 진행하는 방법
- 약한 분류기를 결합하여 강한 분류기를 만드는 과정



1. Gradient Boosting

❖ Boosting이란?

- 부스팅의 대표적인 모델은 AdaBoost, Gradient Boost 등
- Gradient Boost의 변형 모델로는 XGBoost, LightGBM, CatBoost 등

비교	특징
Adaboost	정답분류 및 오답에 가중치 부여
Gradient Boosting	lossFunction에 gradient를 통해 오답에 가중치 부여
Xgboost	GBM 대비 성능향상 시스템 자원 효율적 운용 Kaggle을 통한 성능검증(상위랭커)
Light GBM	Xgboost 대비 성능향상 2000개 이상의 대용량 데이터 사용가능

1. Gradient Boosting

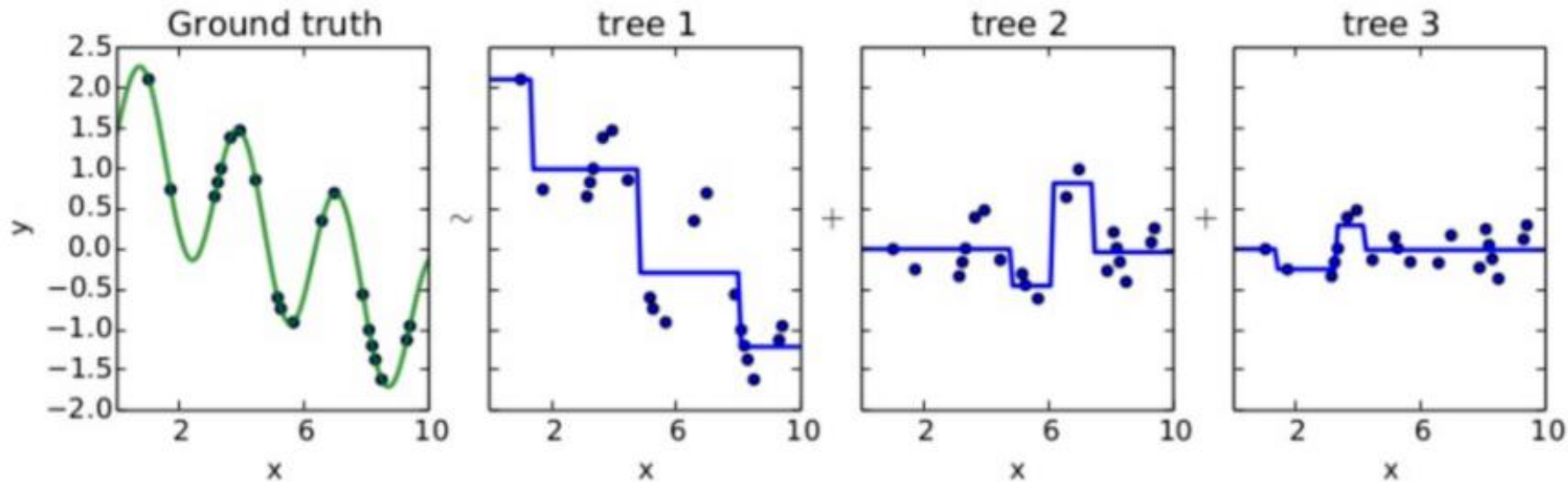
❖ Boosting과 bagging 비교

비교	Bagging	Boosting
특징	병렬 앙상블모델	연속 앙상블 모델
적합한 상황	High Variance , Low Bias	Low variance, High bias
대표알고리즘	RandomForest	Gradient Boosting Xgboosting
Sampling	Random sampling	RandomSampling with on error

1. Gradient Boosting

❖ Gradient Boosting

- 약한 여러 개의 결정트리를 이용하여 강력한 모델을 만듦
- Residual fitting의 방식 예측하는 타겟 추정 값은 모든 타겟 값의 평균
tree1 를 통해 y 를 예측하고 남은 잔차 (residual)을 다시 tree2 라는 모델을 통해 예측하고
tree1+tree2 모델을 통해 y 를 예측한다면 tree1보다 나은 tree2 모델
tree2+tree3 모델을 통해 y 를 예측한다면 tree2보다 나은 tree3 모델



1. Gradient Boosting

❖ Gradient Boosting

- 약한 여러 개의 결정트리를 이용하여 강력한 모델을 만듦
- powerful한 pre-pruning이 사용되며 1~5 정도 깊이의 tree를 사용하므로 메모리를 적게 사용하고 예측도 빠름
- parameter설정에 random forest보다 조금 더 민감하지만 잘 조정하면 높은 정확도를 제공
- parameter는 이전 트리의 오차를 얼마나 강하게 보정할 것인가를 제어하는 **learning_rate**

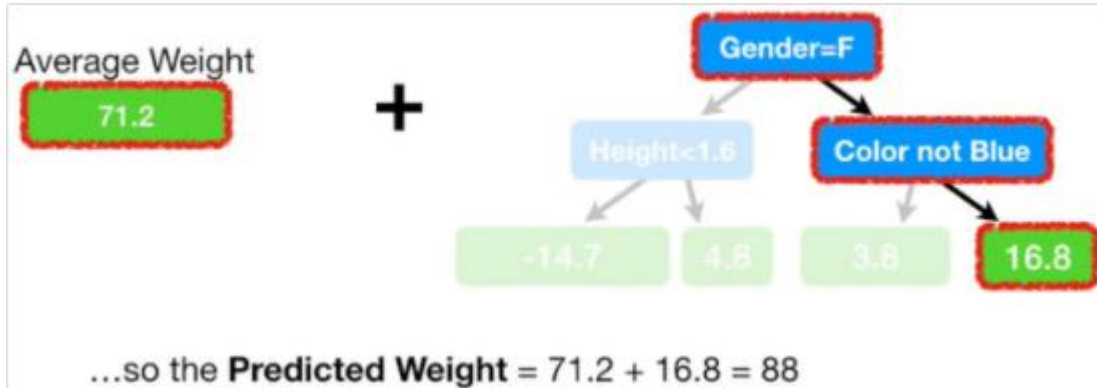
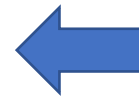
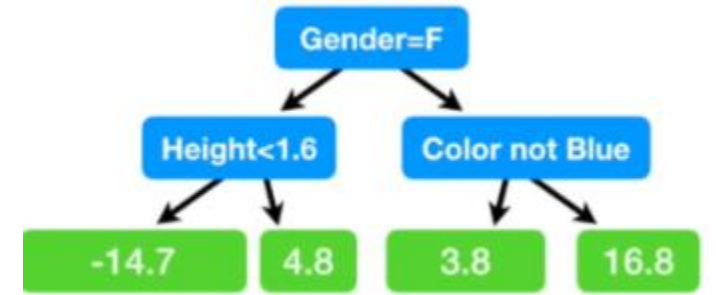
1. Gradient Boosting

❖ Gradient Boosting

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
1.5	Blue	Female	56
1.8	Red	Male	73
1.5	Green	Male	77
1.4	Blue	Female	57



Height (m)	Favorite Color	Gender	Weight (kg)	Residual
1.6	Blue	Male	88	16.8
1.6	Green	Female	76	4.8
1.5	Blue	Female	56	-15.2
1.8	Red	Male	73	1.8
1.5	Green	Male	77	5.8
1.4	Blue	Female	57	-14.2



초기에 구한 트리(여기서는 single leaf)와 두번째로 구한 트리를 조합하여 몸무게를 예측

1. Gradient Boosting

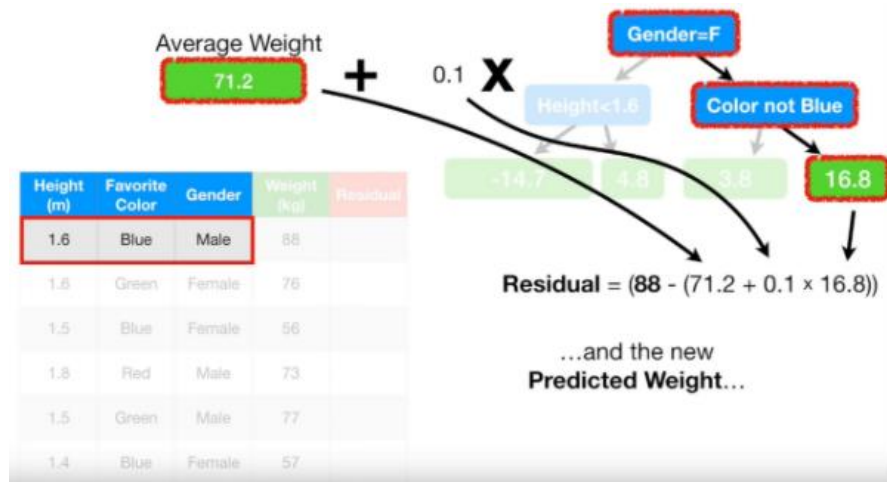
❖ 과적합 문제를 해결하기 위해 학습률(Learning Rate)이라는 것을 활용



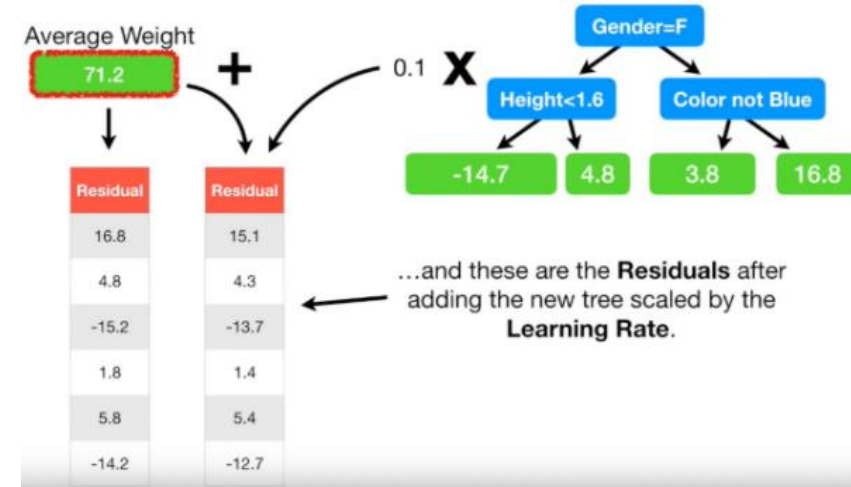
0.1일 경우 모델이 예측한 몸무게는 72.9
single leaf 모델이 예측한 71.2kg보다는 실제 값에 근접

1. Gradient Boosting

❖ 과적합 문제를 해결하기 위해 학습률(Learning Rate)이라는 것을 활용

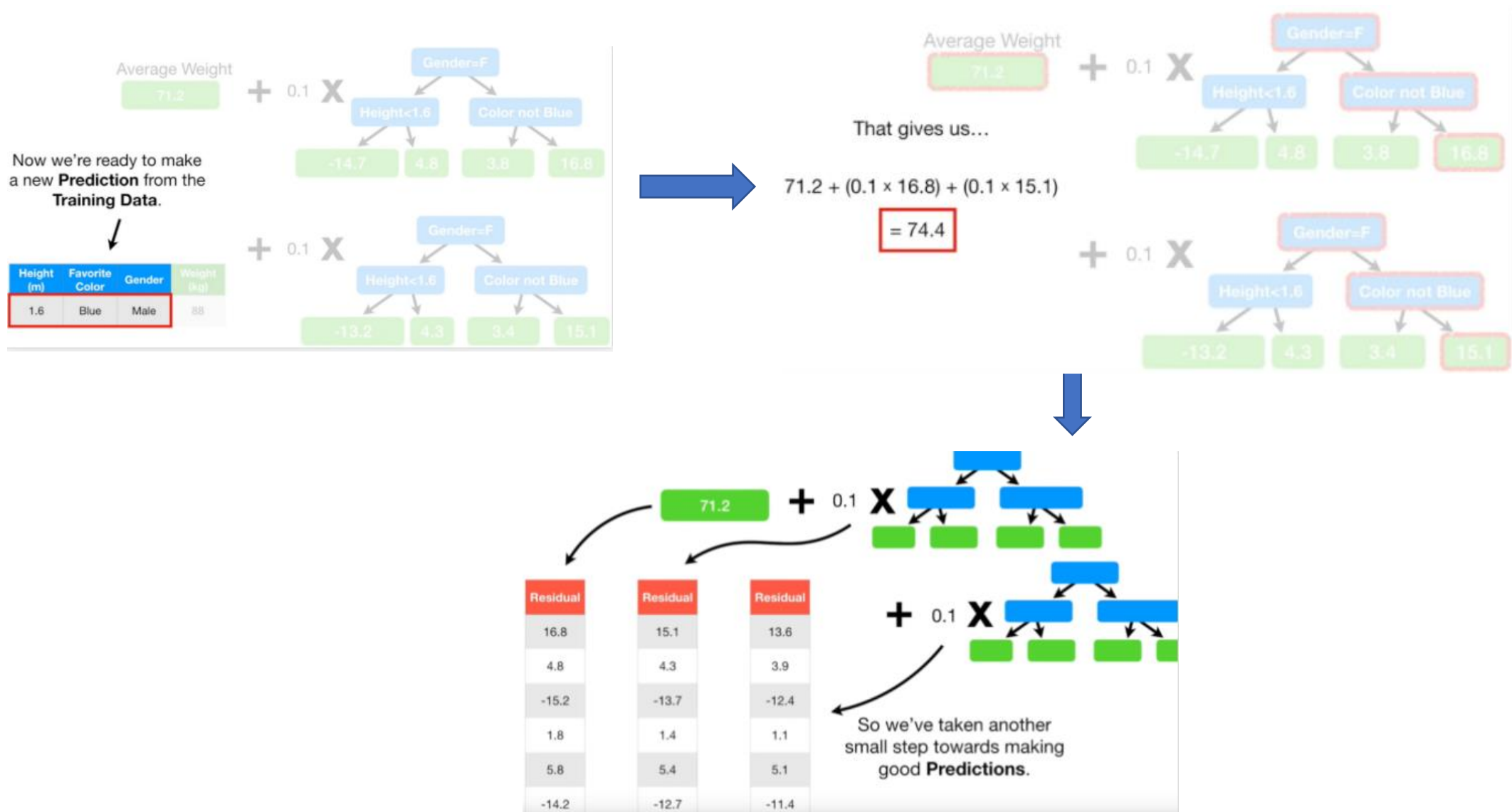


Height (m)	Favorite Color	Gender	Weight (kg)	Residual
1.6	Blue	Male	88	15.1
1.6	Green	Female	76	4.3
1.5	Blue	Female	56	-13.7
1.8	Red	Male	73	1.4
1.5	Green	Male	77	5.4
1.4	Blue	Female	57	-12.7



1. Gradient Boosting

❖ 과적합 문제를 해결하기 위해 학습률(Learning Rate)이라는 것을 활용



1. Gradient Boosting

❖ Gradient Boosting 실습

- `from sklearn.datasets import make_classification`
- `from sklearn.ensemble import GradientBoostingClassifier`
- `clf = GradientBoostingClassifier(random_state=0)`
- `clf.fit(X_train, y_train)`
- `clf.predict(X_test[:2])`
- `clf.score(X_test, y_test)`
-

1. Gradient Boosting

❖ Gradient Boosting 실습

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.datasets import load_breast_cancer
3 import pandas as pd
4 import numpy as np
5 from sklearn.metrics import accuracy_score
6 import matplotlib.pyplot as plt
7
8 cancer = load_breast_cancer()
9 col_names = cancer.feature_names
10 print(len(col_names))
11 X_Data = pd.DataFrame(cancer.data, columns = col_names )
12 y = pd.DataFrame(cancer.target)
13
14 X_train, X_test, y_train, y_test = train_test_split(
15     cancer.data, cancer.target, test_size=0.3, random_state=0)
```

30

1. Gradient Boosting

❖ Gradient Boosting 실습

```
1 from sklearn.ensemble import AdaBoostClassifier
2 ada = AdaBoostClassifier(n_estimators=100, random_state=42)
3 ada.fit(X_train, y_train)
```

AdaBoostClassifier(n_estimators=100, random_state=42)

```
1 print("훈련 세트 정확도: {:.3f}".format(ada.score(X_train, y_train)))
2 print("테스트 세트 정확도: {:.3f}".format(ada.score(X_test, y_test)))
```

훈련 세트 정확도: 1.000
테스트 세트 정확도: 0.977

1. Gradient Boosting

❖ Gradient Boosting 실습

```
1 feature_imp = pd.Series(ada.feature_importances_, index=cancer.feature_names).sort_values(ascending=True)
2 feature_imp
3
4 feature_imp.plot(kind='barh', grid=True, figsize=(5,5))
5 plt.show()
6
```



1. Gradient Boosting

❖ Gradient Boosting 실습

```
1 from sklearn.ensemble import GradientBoostingClassifier
2 gbrt = GradientBoostingClassifier(random_state=0)
3 gbrt.fit(X_train, y_train)
4
5 print("훈련 세트 정확도: {:.3f}".format(gbrt.score(X_train, y_train)))
6 print("테스트 세트 정확도: {:.3f}".format(gbrt.score(X_test, y_test)))
```

훈련 세트 정확도: 1.000
테스트 세트 정확도: 0.977

```
1 gbrt = GradientBoostingClassifier(random_state=0, max_depth=2)
2 gbrt.fit(X_train, y_train)
3
4 print("훈련 세트 정확도: {:.3f}".format(gbrt.score(X_train, y_train)))
5 print("테스트 세트 정확도: {:.3f}".format(gbrt.score(X_test, y_test)))
```

훈련 세트 정확도: 1.000
테스트 세트 정확도: 0.965

```
1 gbrt = GradientBoostingClassifier(random_state=0, learning_rate=0.01)
2 gbrt.fit(X_train, y_train)
3
4 print("훈련 세트 정확도: {:.3f}".format(gbrt.score(X_train, y_train)))
5 print("테스트 세트 정확도: {:.3f}".format(gbrt.score(X_test, y_test)))
```

훈련 세트 정확도: 0.987
테스트 세트 정확도: 0.959

1. Gradient Boosting

❖ Gradient Boosting 실습

```
1 lr_list = [0.05, 0.075, 0.1, 0.25, 0.5, 0.75, 1]
2
3 for learning_rate in lr_list:
4     gb_clf = GradientBoostingClassifier(n_estimators=50, learning_rate=learning_rate, max_depth=2, random_state=0)
5     gb_clf.fit(X_train, y_train)
6
7     print("Learning rate: ", learning_rate)
8     print("Accuracy score (training): {0:.3f}".format(gb_clf.score(X_train, y_train)))
9     print("Accuracy score (test): {0:.3f}".format(gb_clf.score(X_test, y_test)))
```

Learning rate: 0.05
Accuracy score (training): 0.987
Accuracy score (test): 0.947
Learning rate: 0.075
Accuracy score (training): 0.992
Accuracy score (test): 0.947
Learning rate: 0.1
Accuracy score (training): 0.995
Accuracy score (test): 0.947
Learning rate: 0.25
Accuracy score (training): 1.000
Accuracy score (test): 0.965
Learning rate: 0.5
Accuracy score (training): 1.000
Accuracy score (test): 0.965
Learning rate: 0.75
Accuracy score (training): 1.000
Accuracy score (test): 0.959
Learning rate: 1
Accuracy score (training): 1.000
Accuracy score (test): 0.971

1. Gradient Boosting

❖ Gradient Boosting 실습

```
1 from sklearn.metrics import classification_report
2 gb_clf2 = GradientBoostingClassifier(n_estimators=50, learning_rate=0.25, max_depth=2, random_state=0)
3 gb_clf2.fit(X_train, y_train)
4 predictions = gb_clf2.predict(X_test)
5
6 print("Classification Report")
7 print(classification_report(y_test, predictions))
```

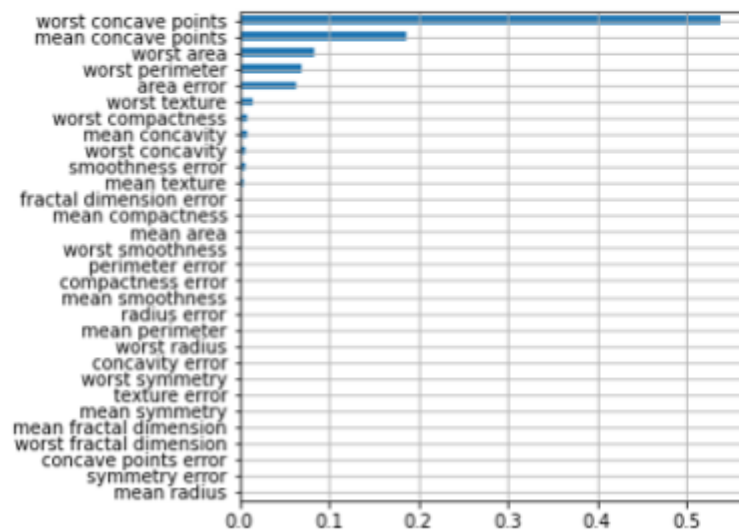
Classification Report

	precision	recall	f1-score	support
0	0.97	0.94	0.95	63
1	0.96	0.98	0.97	108
accuracy			0.96	171
macro avg	0.97	0.96	0.96	171
weighted avg	0.96	0.96	0.96	171

1. Gradient Boosting

❖ Gradient Boosting 실습

```
1 feature_imp = pd.Series(gb_clf2.feature_importances_, index=cancer.feature_names).sort_values(ascending=True)
2 feature_imp
3 # print(feature_imp)
4 feature_imp.plot(kind='barh', grid=True, figsize=(5,5))
5 plt.show()
6
```



1. Gradient Boosting

❖ Gradient Boosting 실습

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.datasets import load_breast_cancer
3 import pandas as pd
4 import numpy as np
5 from sklearn.metrics import accuracy_score
6 import matplotlib.pyplot as plt
7 import mglearn
8 from sklearn.datasets import load_wine
9
10 wine = load_wine()
11 col_names = wine.feature_names
12 print(len(col_names))
13 X_Data = pd.DataFrame(wine.data, columns = col_names )
14 y = pd.DataFrame(wine.target)
15 X_train, X_test, y_train, y_test = train_test_split(
16     wine.data, wine.target, test_size=0.3, random_state=0)
```

13

1. Gradient Boosting

❖ Gradient Boosting 실습

```
1 from sklearn.ensemble import AdaBoostClassifier
2 ada = AdaBoostClassifier(n_estimators=100, random_state=42)
3 ada.fit(X_train, y_train)
```

AdaBoostClassifier(n_estimators=100, random_state=42)

```
1 print("훈련 세트 정확도: {:.3f}".format(ada.score(X_train, y_train)))
2 print("테스트 세트 정확도: {:.3f}".format(ada.score(X_test, y_test)))
```

훈련 세트 정확도: 0.976

테스트 세트 정확도: 0.889

```
1 ada.fit(X_train, y_train)
2 predictions = ada.predict(X_test)
3
4 print("Confusion Matrix:")
5 print(confusion_matrix(y_test, predictions))
6
7 print("Classification Report")
8 print(classification_report(y_test, predictions))
```

Confusion Matrix:

```
[[18  1  0]
 [ 2 18  2]
 [ 0  1 12]]
```

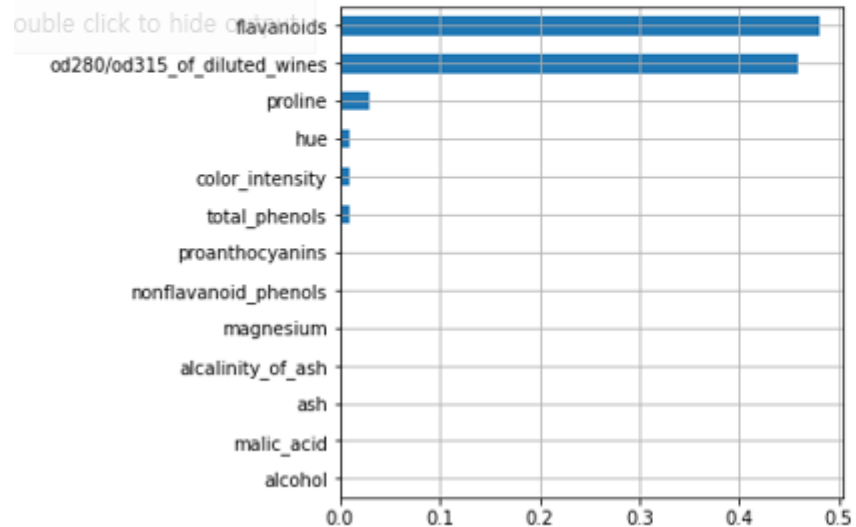
Classification Report

	precision	recall	f1-score	support
0	0.90	0.95	0.92	19
1	0.90	0.82	0.86	22
2	0.86	0.92	0.89	13
accuracy			0.89	54
macro avg	0.89	0.90	0.89	54
weighted avg	0.89	0.89	0.89	54

1. Gradient Boosting

❖ Gradient Boosting 실습

```
1 feature_imp = pd.Series(ada.feature_importances_, index=wine.feature_names).sort_values(ascending=True)
2 feature_imp
3
4 feature_imp.plot(kind='barh', grid=True, figsize=(5,5))
5 plt.show()
6
```



1. Gradient Boosting

❖ Gradient Boosting 실습

```
1 from sklearn.ensemble import GradientBoostingClassifier
2 gbrt = GradientBoostingClassifier(random_state=0)
3 gbrt.fit(X_train, y_train)
4
5 print("훈련 세트 정확도: {:.3f}".format(gbrt.score(X_train, y_train)))
6 print("테스트 세트 정확도: {:.3f}".format(gbrt.score(X_test, y_test)))
```

훈련 세트 정확도: 1.000
테스트 세트 정확도: 0.963

```
1 gbrt = GradientBoostingClassifier(random_state=0, learning_rate=0.01)
2 gbrt.fit(X_train, y_train)
3
4 print("훈련 세트 정확도: {:.3f}".format(gbrt.score(X_train, y_train)))
5 print("테스트 세트 정확도: {:.3f}".format(gbrt.score(X_test, y_test)))
```

훈련 세트 정확도: 1.000
테스트 세트 정확도: 0.944

1. Gradient Boosting

❖ Gradient Boosting 실습

```
1 lr_list = [0.05, 0.075, 0.1, 0.25, 0.5, 0.75, 1]
2
3 for learning_rate in lr_list:
4     gb_clf = GradientBoostingClassifier(n_estimators=50, learning_rate=learning_rate, max_depth=2, random_state=0)
5     gb_clf.fit(X_train, y_train)
6
7     print("Learning rate: ", learning_rate)
8     print("Accuracy score (training): {0:.3f}".format(gb_clf.score(X_train, y_train)))
9     print("Accuracy score (test): {0:.3f}".format(gb_clf.score(X_test, y_test)))
```

```
Learning rate: 0.05
Accuracy score (training): 1.000
Accuracy score (test): 0.963
Learning rate: 0.075
Accuracy score (training): 1.000
Accuracy score (test): 0.963
Learning rate: 0.1
Accuracy score (training): 1.000
Accuracy score (test): 0.981
Learning rate: 0.25
Accuracy score (training): 1.000
Accuracy score (test): 0.963
Learning rate: 0.5
Accuracy score (training): 1.000
Accuracy score (test): 0.963
Learning rate: 0.75
Accuracy score (training): 1.000
Accuracy score (test): 0.963
Learning rate: 1
Accuracy score (training): 1.000
Accuracy score (test): 0.963
```

1. Gradient Boosting

❖ Gradient Boosting 실습

```
1 from sklearn.metrics import classification_report
2 gb_clf2 = GradientBoostingClassifier(n_estimators=50, learning_rate=0.25, max_depth=2, random_state=0)
3 gb_clf2.fit(X_train, y_train)
4 predictions = gb_clf2.predict(X_test)
5
6 print("Classification Report")
7 print(classification_report(y_test, predictions))
```

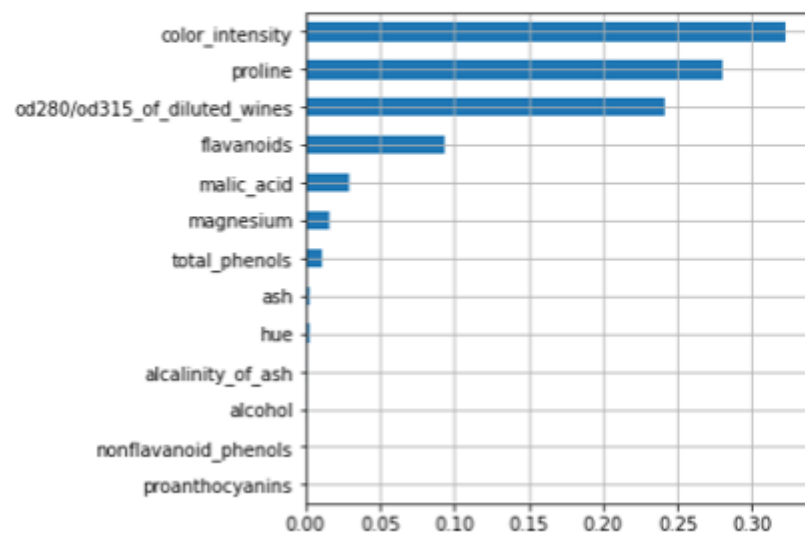
Classification Report

	precision	recall	f1-score	support
0	0.96	1.00	0.97	19
1	1.00	0.91	0.95	22
2	0.93	1.00	0.96	13
accuracy			0.96	54
macro avg	0.96	0.97	0.96	54
weighted avg	0.97	0.96	0.96	54

1. Gradient Boosting

❖ Gradient Boosting 실습

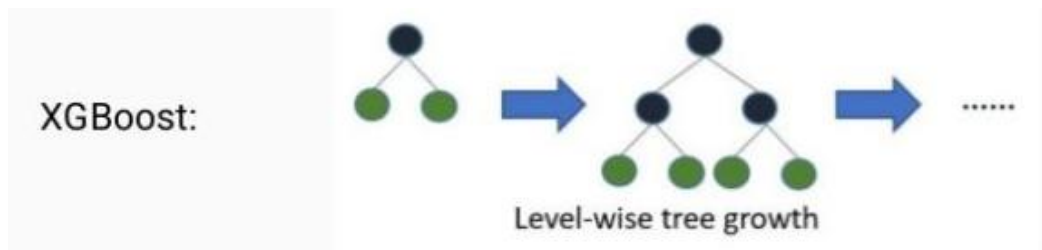
```
1 feature_imp = pd.Series(gb_clf2.feature_importances_, index=wine.feature_names).sort_values(ascending=True)
2 feature_imp
3 # print(feature_imp)
4 feature_imp.plot(kind='barh', grid=True, figsize=(5,5))
5 plt.show()
6
```



1. Gradient Boosting

❖ xgBoosting

- 트리기반의 앙상블 학습중 가장 최고의 모델로 **Gradient Boosting**보다 빠름
- 과적합을 방지와 분류. 회귀 둘다 가능하고 조기 종료(early stopping) 를 제공함
- **Gradient boosting** 기반으로 학습시간이 느림/ 하이퍼 파라미터 존재
- Level-wise : 각 노드는 root노드와 가까운노드를 선회 / 수평성장
- 최대한 균형 잡힌 트리를 유지하며 분할하여 트리의 깊이를 최소화하여 오버피팅에 강한구조이지만 균형을 맞추기 위한 시간필요



`n_estimators`(혹은 `num_boost_round`) : 결정 트리의 개수

`max_depth` : 트리의 깊이

`colsample_bytree` : 컬럼의 샘플링 비율(`random forest`의 `max_features`와 비슷)

`subsample` : weak learner가 학습에 사용하는 데이터 샘플링 비율

`learning_rate` : 학습률

`min_split_loss` : 리프 노드를 추가적으로 나눌지 결정하는 값

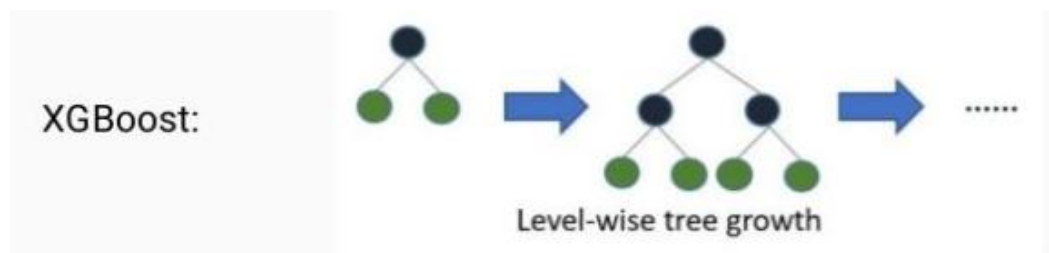
`reg_lambda` : L2 규제

`reg_alpha` : L1 규제

1. Gradient Boosting

❖ xgBoosting

- GBM의 경우 `n_estimators`에 지정된 횟수만큼 학습을 끝까지 수행하지만, XGB의 경우 오류가 더 이상 개선되지 않으면 수행을 중지
- `n_estimators` 를 200으로 설정하고, 조기 중단 파라미터 값을 50으로 설정하면, 1부터 200회까지 부스팅을 반복후 학습오류가 감소하지 않으면 더 이상 부스팅을 진행하지 않고 종료함



`n_estimators`(혹은 `num_boost_round`) : 결정 트리의 개수
`max_depth` : 트리의 깊이
`colsample_bytree` : 컬럼의 샘플링 비율(`random forest`의 `max_features`와 비슷)
`subsample` : weak learner가 학습에 사용하는 데이터 샘플링 비율
`learning_rate` : 학습률
`min_split_loss` : 리프 노드를 추가적으로 나눌지 결정하는 값
`reg_lambda` : L2 규제
`reg_alpha` : L1 규제

1. Gradient Boosting

❖ xgBoosting 설치

- #콘다환경에서 설치
- `conda install -c anaconda py-xgboost`
- `pip install xgboost`

1. Gradient Boosting

❖ xgBoosting 실습

- `from xgboost import plot_importance`
- `from xgboost import XGBClassifier`
- `clf=xgb.XGBClassifier()` # 파라미터 넣어줌. 모델생성
- `clf.fit()` # 파라미터 넣어줌. 데이터학습.
- `evals_result=clf.evals_result()`
- `clf.feature_importances_`
- `clf.predict()`

1. Gradient Boosting

❖ xgBoosting 실습

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.datasets import load_breast_cancer
3 import pandas as pd
4 import numpy as np
5 from sklearn.metrics import accuracy_score
6 import matplotlib.pyplot as plt
7
8 cancer = load_breast_cancer()
9 col_names = cancer.feature_names
10 print(len(col_names))
11 X_Data = pd.DataFrame(cancer.data, columns = col_names )
12 y = pd.DataFrame(cancer.target)
13
14 X_train, X_test, y_train, y_test = train_test_split(
15     cancer.data, cancer.target, test_size=0.3, random_state=0)
```

30

1. Gradient Boosting

❖ xgBoosting 실습

```
1 from xgboost import plot_importance
2 from xgboost import XGBClassifier
3
4 # import xgboost as xgb
```

```
1 xgbb = XGBClassifier(n_estimators = 400, learning_rate = 0.1, max_depth = 3)
2 X_train = pd.DataFrame(X_train, columns=cancer.feature_names)
3 X_test = pd.DataFrame(X_test, columns=cancer.feature_names)
4 xgbb.fit(X_train, y_train)
```

XGBClassifier(n_estimators=400)

```
1 print("훈련 세트 정확도: {:.3f}".format(xgbb.score(X_train, y_train)))
2 print("테스트 세트 정확도: {:.3f}".format(xgbb.score(X_test, y_test)))
```

훈련 세트 정확도: 1.000
테스트 세트 정확도: 0.982

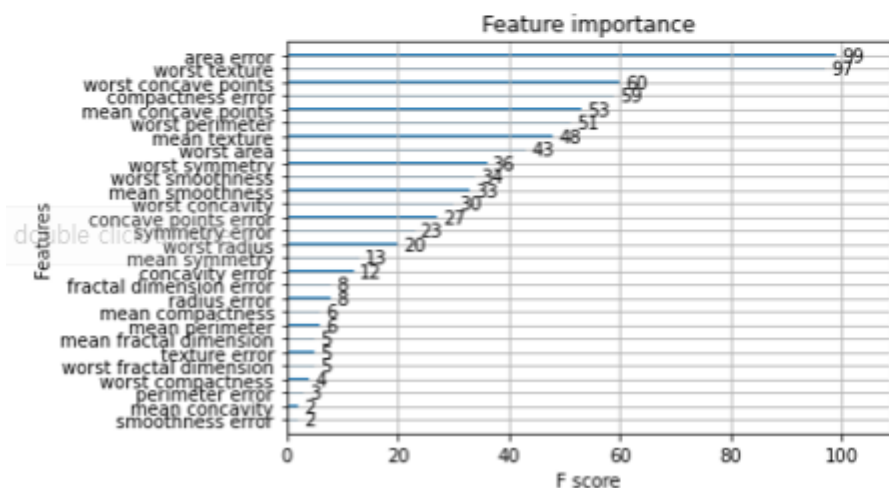
1. Gradient Boosting

❖ xgBoosting 실습

```
1 fscore = xgbb.get_booster().get_fscore()
2 score_f = sorted(fscore.items(), key=(lambda x: x[1]), reverse=True)
3 print(score_f)
4
```

```
[('area error', 99), ('worst texture', 97), ('worst concave points', 60), ('compactness error', 59), ('mean concave points', 53), ('worst p
erimeter', 51), ('mean texture', 48), ('worst area', 43), ('worst symmetry', 36), ('worst smoothness', 34), ('mean smoothness', 33), ('wors
t concavity', 30), ('concave points error', 27), ('symmetry error', 23), ('worst radius', 20), ('mean symmetry', 13), ('concavity error', 1
2), ('radius error', 8), ('fractal dimension error', 8), ('mean perimeter', 6), ('mean compactness', 6), ('worst fractal dimension', 5),
('texture error', 5), ('mean fractal dimension', 5), ('worst compactness', 4), ('perimeter error', 3), ('smoothness error', 2), ('mean conc
avity', 2)]
```

```
1 import matplotlib.pyplot as plt
2 plot_importance(xgbb)
3 plt.show()
```



1. Gradient Boosting

❖ xgBoosting 실습

```
1 evals= [(X_test, y_test)]

1 xgbb.fit(X_train, y_train, early_stopping_rounds=100, eval_metric= "logloss", eval_set=evals, verbose=True)
[332]    validation_0-logloss:0.067166
[333]    validation_0-logloss:0.067109
[334]    validation_0-logloss:0.066941
[335]    validation_0-logloss:0.067048
[336]    validation_0-logloss:0.066874
[337]    validation_0-logloss:0.066998
[338]    validation_0-logloss:0.066829
[339]    validation_0-logloss:0.066813
[340]    validation_0-logloss:0.066918
[341]    validation_0-logloss:0.066824
[342]    validation_0-logloss:0.066766
[343]    validation_0-logloss:0.066889
[344]    validation_0-logloss:0.066725
[345]    validation_0-logloss:0.066832
[346]    validation_0-logloss:0.06666
[347]    validation_0-logloss:0.066604
[348]    validation_0-logloss:0.066515
[349]    validation_0-logloss:0.066619
```

```
XGBClassifier(n_estimators=400)
```


1. Gradient Boosting

❖ xgBoosting 실습

```
1 from sklearn.metrics import confusion_matrix, accuracy_score, precision_score, recall_score
2 from sklearn.metrics import confusion_matrix, f1_score, roc_auc_score
3 def get_clf_eval(y_test, y_pred):
4     accuracy = accuracy_score(y_test, y_pred)
5     precision = precision_score(y_test, y_pred)
6     recall = recall_score(y_test, y_pred)
7     F1 = f1_score(y_test, y_pred)
8     print('ㄱ정확도: {:.4f}'.format(accuracy))
9     print('정밀도: {:.4f}'.format(precision))
10    print('재현율: {:.4f}'.format(recall))
11    print('F1: {:.4f}'.format(F1))
12
```

```
1
2 xgb_wrapper = XGBClassifier(learning_rate=0.01, max_depth=5, n_estimators=400,
3                             random_state=42)
4
5 evals = [(X_test, y_test)]
6 xgb_wrapper.fit(X_train, y_train, early_stopping_rounds=10,
7                eval_metric='logloss', eval_set=evals, verbose=True)
8
9 wsl0_preds = xgb_wrapper.predict(X_test)
10 get_clf_eval(y_test, wsl0_preds)
```

```
[386] validation_0-logloss:0.100167
[387] validation_0-logloss:0.100036
[388] validation_0-logloss:0.099796
[389] validation_0-logloss:0.099652
[390] validation_0-logloss:0.09954
[391] validation_0-logloss:0.099366
[392] validation_0-logloss:0.099132
[393] validation_0-logloss:0.099062
[394] validation_0-logloss:0.098937
[395] validation_0-logloss:0.098708
[396] validation_0-logloss:0.098572
[397] validation_0-logloss:0.098464
[398] validation_0-logloss:0.098213
[399] validation_0-logloss:0.098149
```

정확도: 0.9649
정밀도: 0.9722
재현율: 0.9722
F1: 0.9722

1. Gradient Boosting

❖ xgBoosting 실습

```
1 from sklearn.model_selection import GridSearchCV
```

```
1 xg = XGBClassifier()  
2  
3 param_grid={  
4     'max_depth': [4, 6, 8, 10, 12],  
5     'n_estimators': [50, 100],  
6     'learning_rate': [0.01, 0.05, 0.1, 0.15]}
```

```
1 gcv=GridSearchCV(xg, param_grid=param_grid, cv=5, scoring='accuracy', n_jobs=1)  
2 gcv.fit(X_train, y_train)
```

```
GridSearchCV(cv=5, estimator=XGBClassifier(), n_jobs=1,  
             param_grid={'learning_rate': [0.01, 0.05, 0.1, 0.15],  
                         'max_depth': [4, 6, 8, 10, 12],  
                         'n_estimators': [50, 100]},  
             scoring='accuracy')
```

```
1 print('final params', gcv.best_params_) # 최적의 파라미터 값 출력  
2 print('best score', gcv.best_score_)
```

```
final params {'learning_rate': 0.01, 'max_depth': 6, 'n_estimators': 100}  
best score 0.949746835443038
```

1. Gradient Boosting

❖ xgBoosting 실습

```
1 cv_result_df=pd.DataFrame(gcv.cv_results_)
2 cv_result_df.sort_values(by=['rank_test_score'], inplace=True)
3
```

```
1 cv_result_df[['params', 'mean_test_score', 'rank_test_score']].head(10)
```

	params	mean_test_score	rank_test_score
19	{'learning_rate': 0.05, 'max_depth': 12, 'n_est...	0.949747	1
37	{'learning_rate': 0.15, 'max_depth': 10, 'n_est...	0.949747	1
35	{'learning_rate': 0.15, 'max_depth': 8, 'n_esti...	0.949747	1
33	{'learning_rate': 0.15, 'max_depth': 6, 'n_esti...	0.949747	1
29	{'learning_rate': 0.1, 'max_depth': 12, 'n_esti...	0.949747	1
27	{'learning_rate': 0.1, 'max_depth': 10, 'n_esti...	0.949747	1
25	{'learning_rate': 0.1, 'max_depth': 8, 'n_estim...	0.949747	1
23	{'learning_rate': 0.1, 'max_depth': 6, 'n_estim...	0.949747	1
17	{'learning_rate': 0.05, 'max_depth': 10, 'n_est...	0.949747	1
15	{'learning_rate': 0.05, 'max_depth': 8, 'n_esti...	0.949747	1

1. Gradient Boosting

❖ xgBoosting 실습

```
: 1 xg_cv = XGBClassifier(learning_rate=0.01, max_depth= 6, n_estimators= 100)
```

```
: 1 xg_cv.fit(X_train, y_train)
```

```
: XGBClassifier(learning_rate=0.01, max_depth=6)
```

```
: 1 predictions = xg_cv.predict(X_test)
  2
  3 print("Classification Report")
  4 print(classification_report(y_test, predictions))
```

```
Classification Report
              precision    recall  f1-score   support

     0       0.98         0.95         0.97         63
     1       0.97         0.98         0.98        108

 accuracy          0.98
 macro avg         0.98         0.97         0.97         171
weighted avg         0.98         0.98         0.98         171
```

1. Gradient Boosting

❖ xgBoosting 실습

```
1 from xgboost import XGBRegressor # 회귀트리 모델
2 from xgboost import plot_importance # 중요변수 시각화
3
4 from sklearn.datasets import load_boston # dataset
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import mean_absolute_error, mean_squared_error
```

```
1 boston = load_boston()
2 X = boston.data
3 y = boston.target
4 col_names = boston.feature_names
```

```
1 print(X.shape), # (506, 13) # 13개의 칼럼
2 print(y.shape)
3 print(col_names)
4
```

(506, 13)

(506,)

['CRIM' 'ZN' 'INDUS' 'CHAS' 'NOX' 'RM' 'AGE' 'DIS' 'RAD' 'TAX' 'PTRATIO'
'B' 'LSTAT']

1. Gradient Boosting

❖ xgBoosting 실습

```
1 x_train, x_test, y_train, y_test = train_test_split(X, y, test_size = 0.3)
2 print(x_train.shape) # (354, 13)
3 print(x_test.shape) # (152, 13)
```

(354, 13)

(152, 13)

```
1 model = XGBRegressor(objective='reg:squarederror')
2 model.fit(x_train, y_train)
3
```

XGBRegressor(objective='reg:squarederror')

```
1 y_predd = model.predict(x_test)
2 y_true = y_test
```

```
1 mse = mean_squared_error(y_true, y_predd)
2 mse # 12.75553963355965
```

9.561559682544305

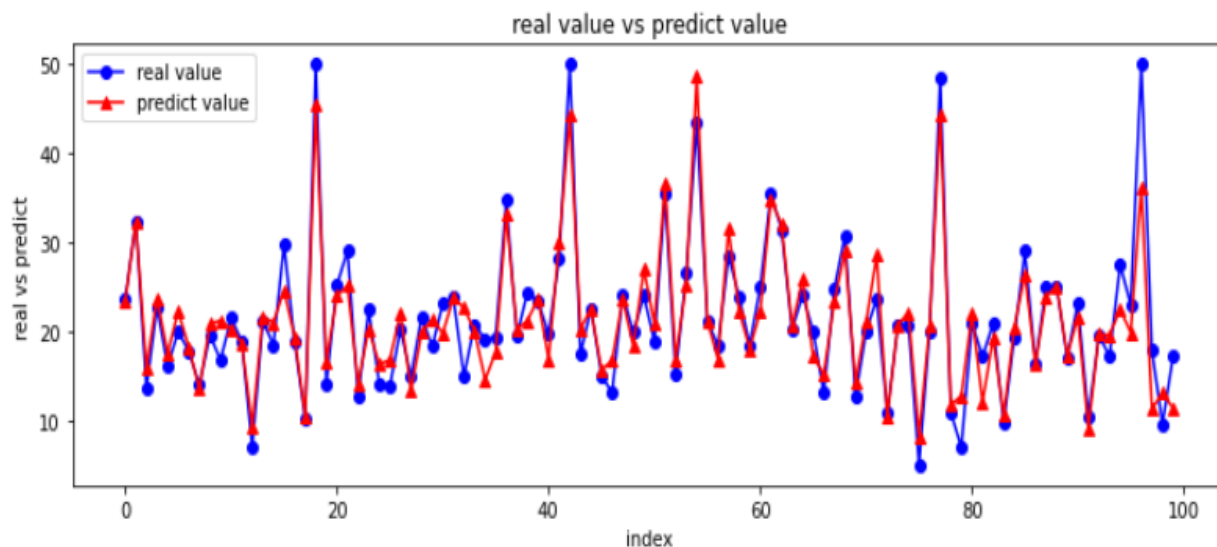
- objective 기본값 : reg:squarederror
- 목적함수이다. 이 함수를 통해 나온 값이 최소화되는 방향으로 학습.
- 종류가 너무 다양해 자주 쓰는 것들만 설명하자면
 - reg:squarederror / reg:squaredlogerror : 오차 제곱 / 오차 로그 제곱
 - binary:logistic : 이항 분류(binary class)에 사용.
 - multi:softmax / multi:softprob : 다항 분류(multi class)에 사용.

1. Gradient Boosting

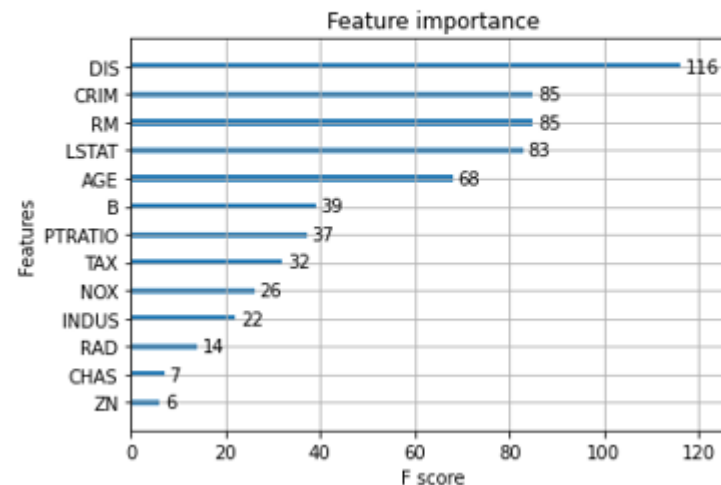
❖ xgBoosting 실습

```
1 import matplotlib.pyplot as plt
2
3 fig = plt.figure( figsize = (12, 4) )
4 chart = fig.add_subplot(1,1,1)
5 chart.plot(y_true[:100], marker='o', color='blue', label='real value')
6 chart.plot(y_predd[:100], marker='^', color='red', label='predict value')
7 chart.set_title('real value vs predict value')
8 plt.xlabel('index')
9 plt.ylabel('real vs predict')
10 plt.legend(loc = 'best')
11
12
```

<matplotlib.legend.Legend at 0x1d09f24fb70>



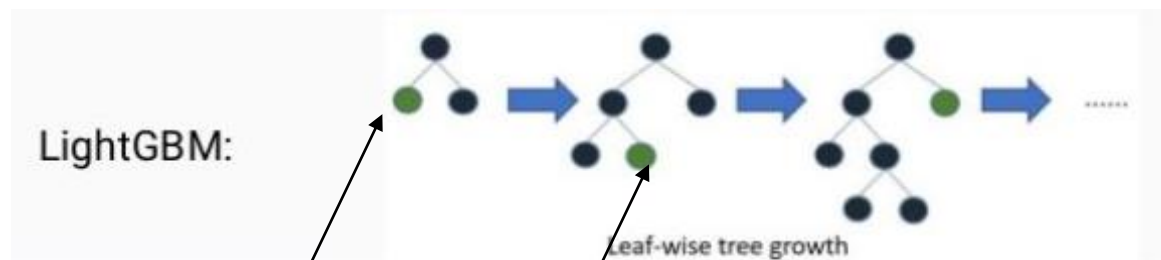
```
1 import matplotlib.pyplot as plt
2 plot_importance(model)
3 plt.show()
```



1. Gradient Boosting

❖ LightGBM

- Xgboosting 을 해결하기 위해 나온 모델로 대용량의 데이터에서 사용 가능
- Level-wise 방법의 Xgboosting은 tree의 균형을 잡아주지만, 연산이 많고 LightGBM은 트리의 균형은 맞추지 않고 리프노드를 지속적으로 분할
- leaf-wise 알고리즘은 level-wise 알고리즘보다 더 많은 loss, 손실을 줄임 Loss 변화가 큰 노드에서 데이터를 분할 성장. 수직 성장
- 학습데이터가 많을 경우 뛰어난 성능



출처 : <https://www.slideshare.net/GabrielCyprianoSaca/xgboost-lightgbm>

Max Loss

Max Loss

n_estimators : 반복하려는 트리의 개수

learning_rate : 학습률

max_depth : 트리의 최대 깊이

min_child_samples : 리프 노드가 되기 위한 최소한의 샘플 데이터 수

num_leaves : 하나의 트리가 가질 수 있는 최대 리프 개수

feature_fraction : 트리를 학습할 때마다 선택하는 feature의 비율

reg_lambda : L2 regularization

reg_alpha : L1 regularization

1. Gradient Boosting

❖ LightGBM

▪ 장점

- XGBoost 대비 더 빠른 학습과 예측 수행 시간
- 더 작은 메모리 사용량
- 카테고리형 피처의 자동 변환과 최적 분할

▪ 단점

- 적은 데이터 세트에 적용할 경우 과적합이 발생하기 쉽고 (공식 문서상 대략 10,000건 이하의 데이터 세트)

1. Gradient Boosting

❖ LightGBM 설치

- #콘다환경에서 설치
- `conda install -c conda-forge lightgbm`
- `pip install lightgbm`

1. Gradient Boosting

❖ LightGBM 실습

- `from lightgbm import LGBMClassifier , plot_importance`
- `model_lgb=LGBMClassifier()`
- `model_lgb.fit(x_train, y_train)`
- `model_lgb. predict(x_train, y_train)`

1. Gradient Boosting

❖ LightGBM 실습

파라미터 명	설명
objective	<ul style="list-style-type: none">- 'reg:linear' : 회귀- binary:logistic : 이진분류- multi:softmax : 다중분류, 클래스 반환- multi:softprob : 다중분류, 확률반환
eval_metric	<ul style="list-style-type: none">- 검증에 사용되는 함수정의- 회귀 분석인 경우 'rmse'를, 클래스 분류 문제인 경우 'error' <p>-----</p> <ul style="list-style-type: none">- rmse : Root Mean Squared Error- mae : mean absolute error- logloss : Negative log-likelihood- error : binary classification error rate- merror : multiclass classification error rate- mlogloss: Multiclass logloss- auc: Area Under Curve
early_stopping_rounds	<p>eval_set : 성능평가를 위한 평가용 데이터 세트를 설정</p> <p>eval_metric : 평가 세트에 적용할 성능 평가 방법 (반복마다 eval_set으로 지정된 데이터 세트에서 eval_metric의 지정된 평가 지표로 예측 오류를 측정)</p>

1. Gradient Boosting

❖ LightGBM

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.datasets import load_breast_cancer
3 import pandas as pd
4 import numpy as np
5 from sklearn.metrics import accuracy_score
6 import matplotlib.pyplot as plt
7
8 cancer = load_breast_cancer()
9 col_names = cancer.feature_names
10 print(len(col_names))
11 X_Data = pd.DataFrame(cancer.data, columns = col_names )
12 y = pd.DataFrame(cancer.target)
13
14 X_train, X_test, y_train, y_test = train_test_split(
15     cancer.data, cancer.target, test_size=0.3, random_state=0)
```

30

1. Gradient Boosting

❖ LightGBM

```
1 from lightgbm import LGBMClassifier, plot_importance
```

```
1 lgb = LGBMClassifier(n_estimators=500, random_state=42)
2 lgb.fit(X_train, y_train)
```

LGBMClassifier(n_estimators=500, random_state=42)

```
1 print("훈련 세트 정확도: {:.3f}".format(lgb.score(X_train, y_train)))
2 print("테스트 세트 정확도: {:.3f}".format(lgb.score(X_test, y_test)))
```

훈련 세트 정확도: 1.000

테스트 세트 정확도: 0.982

```
1 lgb.fit(X_train, y_train)
2 predictions = lgb.predict(X_test)
3
4 print("Classification Report")
5 print(classification_report(y_test, predictions))
```

Classification Report

	precision	recall	f1-score	support
0	1.00	0.95	0.98	63
1	0.97	1.00	0.99	108
accuracy			0.98	171
macro avg	0.98	0.98	0.98	171
weighted avg	0.98	0.98	0.98	171

1. Gradient Boosting

❖ LightGBM

```
1 evals= [(X_test, y_test)]
```

```
1 lgb.fit(X_train, y_train, early_stopping_rounds=100, eval_metric="logloss", eval_set=evals, verbose=True)
```

```
[485] valid_0's binary_logloss: 0.0301314
[486] valid_0's binary_logloss: 0.0301314
[487] valid_0's binary_logloss: 0.0301314
[488] valid_0's binary_logloss: 0.0301314
[489] valid_0's binary_logloss: 0.0301314
[490] valid_0's binary_logloss: 0.0301314
[491] valid_0's binary_logloss: 0.0301314
[492] valid_0's binary_logloss: 0.0301314
[493] valid_0's binary_logloss: 0.0301314
[494] valid_0's binary_logloss: 0.0301314
[495] valid_0's binary_logloss: 0.0301314
[496] valid_0's binary_logloss: 0.0301314
[497] valid_0's binary_logloss: 0.0301314
[498] valid_0's binary_logloss: 0.0301314
[499] valid_0's binary_logloss: 0.0301314
[500] valid_0's binary_logloss: 0.0301314
Did not meet early stopping. Best iteration is:
[407] valid_0's binary_logloss: 0.0300678
```

```
LGBMClassifier(n_estimators=500, random_state=42)
```

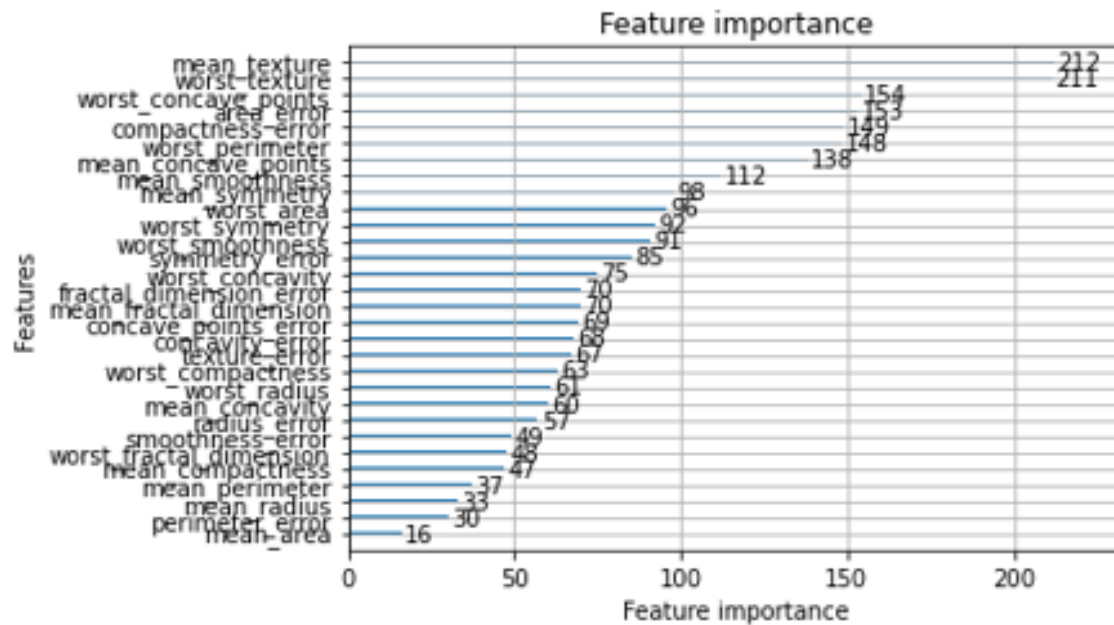
early_stopping_rounds 파라미터 : 조기 중단을 위한 라운드를 설정합니다.
조기 중단 기능 수행을 위해서는 반드시 **eval_set**과 **eval_metric**이 함께 설정되어야 합니다.

- **eval_set** : 성능평가를 위한 평가용 데이터 세트를 설정
- **eval_metric** : 평가 세트에 적용할 성능 평가 방법
(반복마다 **eval_set**으로 지정된 데이터 세트에서 **eval_metric**의 지정된 평가 지표로 예측 오류를 측정)

1. Gradient Boosting

❖ LightGBM

```
1 import matplotlib.pyplot as plt
2 plot_importance(lgb)
3 plt.show()
```



1. Gradient Boosting

❖ LightGBM

```
1 lb = LGBMClassifier()
2
3 param_grid={
4     'max_depth': [1,2,3,4],
5     'n_estimators': [50,100],
6     'learning_rate': [0.01,0.05,0.1,0.15]}
```

```
1 lgb_cv=GridSearchCV(lb, param_grid=param_grid, cv=5, scoring='accuracy', n_jobs=1)
2 lgb_cv.fit(X_train, y_train)
```

```
GridSearchCV(cv=5, estimator=LGBMClassifier(), n_jobs=1,
             param_grid={'learning_rate': [0.01, 0.05, 0.1, 0.15],
                          'max_depth': [1, 2, 3, 4], 'n_estimators': [50, 100]},
             scoring='accuracy')
```

```
1 print('final params', lgb_cv.best_params_) # 최적의 파라미터 값 출력
2 print('best score', lgb_cv.best_score_)
```

```
final params {'learning_rate': 0.15, 'max_depth': 2, 'n_estimators': 100}
best score 0.967373417721519
```

1. Gradient Boosting

❖ LightGBM

```
1 cv_result_df=pd.DataFrame(lgb_cv.cv_results_)
2 cv_result_df.sort_values(by=['rank_test_score'], inplace=True)
```

```
1 cv_result_df[['params', 'mean_test_score', 'rank_test_score']].head(10)
```

	params	mean_test_score	rank_test_score
27	{'learning_rate': 0.15, 'max_depth': 2, 'n_est...	0.967373	1
23	{'learning_rate': 0.1, 'max_depth': 4, 'n_esti...	0.964873	2
21	{'learning_rate': 0.1, 'max_depth': 3, 'n_esti...	0.964873	2
28	{'learning_rate': 0.15, 'max_depth': 3, 'n_est...	0.964842	4
29	{'learning_rate': 0.15, 'max_depth': 3, 'n_est...	0.962373	5
25	{'learning_rate': 0.15, 'max_depth': 1, 'n_est...	0.962342	6
31	{'learning_rate': 0.15, 'max_depth': 4, 'n_est...	0.959873	7
30	{'learning_rate': 0.15, 'max_depth': 4, 'n_est...	0.959842	8
19	{'learning_rate': 0.1, 'max_depth': 2, 'n_esti...	0.959778	9
26	{'learning_rate': 0.15, 'max_depth': 2, 'n_est...	0.957310	10

1. Gradient Boosting

❖ LightGBM

```
: 1 lb_cv = LGBMClassifier(learning_rate= 0.15, max_depth=2, n_estimators= 100)
: 2 lb_cv.fit(X_train, y_train)
```

```
: LGBMClassifier(learning_rate=0.15, max_depth=2)
```

```
: 1 predictions = lb_cv.predict(X_test)
: 2
: 3 print("Classification Report")
: 4 print(classification_report(y_test, predictions))
```

Classification Report

	precision	recall	f1-score	support
0	0.97	0.95	0.96	63
1	0.97	0.98	0.98	108
accuracy			0.97	171
macro avg	0.97	0.97	0.97	171
weighted avg	0.97	0.97	0.97	171

1. Gradient Boosting

❖ LightGBM

```
: 1 from lightgbm import LGBMRegressor, plot_importance
```

```
: 1
```

```
: 1 model_lgb=LGBMRegressor()  
2 x_train = pd.DataFrame(x_train, columns=boston.feature_names)  
3 x_test = pd.DataFrame(x_test, columns=boston.feature_names)  
4  
5 model_lgb.fit(x_train, y_train)  
6 model_lgb
```

```
: LGBMRegressor()
```

```
: 1 y_pred = model_lgb.predict(x_test)  
2 y_true = y_test
```

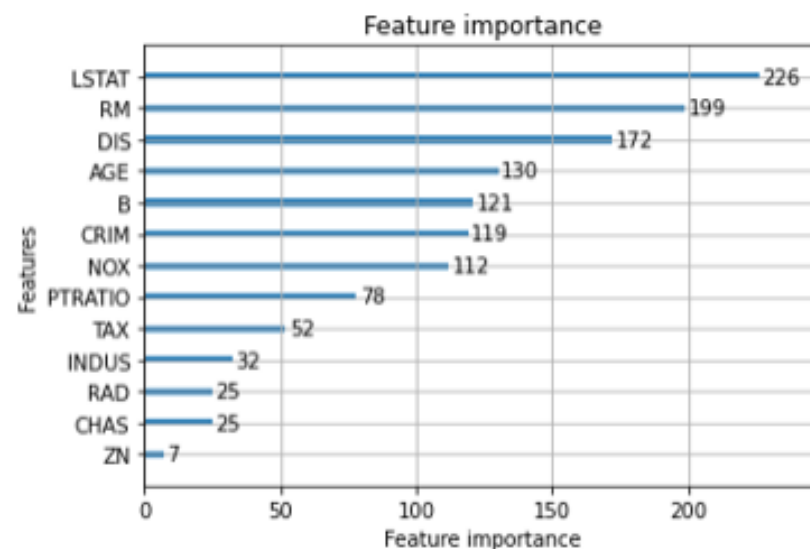
```
: 1 mse = mean_squared_error(y_true, y_pred)  
2 mse
```

```
: 11.346167217013914
```

```
: 1 print(boston.feature_names)
```

```
['CRIM' 'ZN' 'INDUS' 'CHAS' 'NOX' 'RM' 'AGE' 'DIS' 'RAD' 'TAX' 'PTRATIO'  
 'B' 'LSTAT']
```

```
1 import matplotlib.pyplot as plt  
2 plot_importance(model_lgb)  
3 plt.show()
```

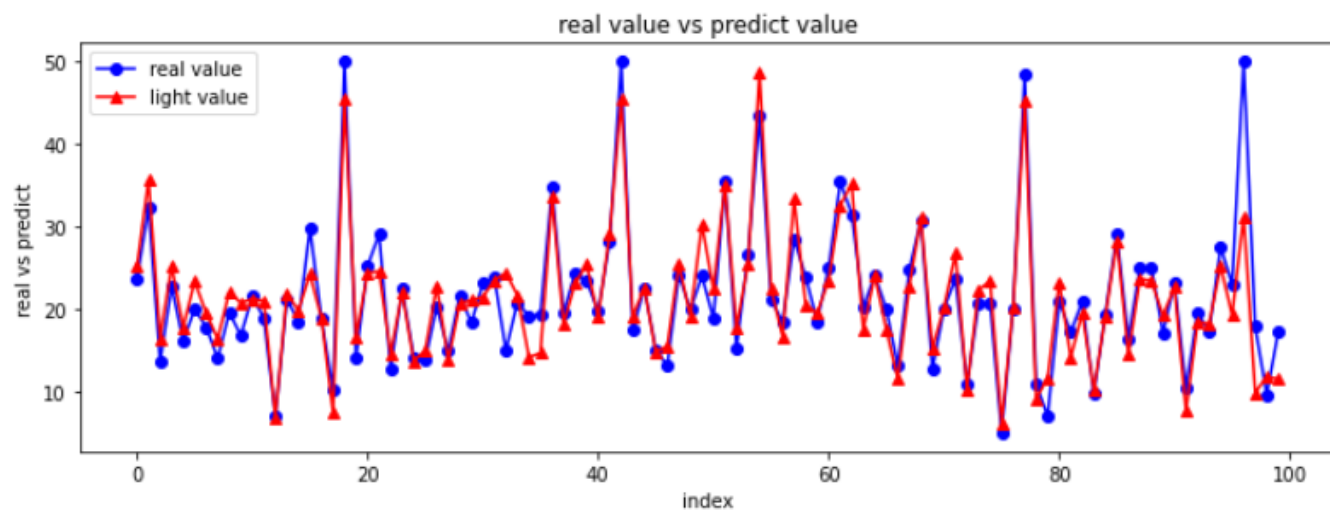


1. Gradient Boosting

❖ LightGBM

```
1 import matplotlib.pyplot as plt
2
3 # y_true.shape # (5160,) : 5160개의 데이터 => 많으니까 100개만 출력 시도
4
5 fig = plt.figure( figsize = (12, 4) )
6 chart = fig.add_subplot(1,1,1)
7 chart.plot(y_true[:100], marker='o', color='blue', label='real value')
8 # chart.plot(y_pred[:100], marker='s', color='green', label='xgboost value')
9 chart.plot(y_pred[:100], marker='^', color='red', label='light value')
10 chart.set_title('real value vs predict value')
11 plt.xlabel('index')
12 plt.ylabel('real vs predict')
13 plt.legend(loc = 'best')
14
15 # 위에 블럭 실행
```

<matplotlib.legend.Legend at 0x11c1fbeb6a0>



1. Gradient Boosting

❖ LightGBM

```
1 model = XGBRegressor(objective='reg:squarederror')
2 model.fit(x_train, y_train)
3
```

XGBRegressor(objective='reg:squarederror')

```
1 y_predd = model.predict(x_test)
2 y_true = y_test
```

```
1 mse = mean_squared_error(y_true, y_predd)
2 mse
```

9.561559682544305

```
1 model_lgb=LGBMRegressor()
2 x_train = pd.DataFrame(x_train, columns=boston.feature_names)
3 x_test = pd.DataFrame(x_test, columns=boston.feature_names)
4
5 model_lgb.fit(x_train, y_train)
6 model_lgb
```

LGBMRegressor()

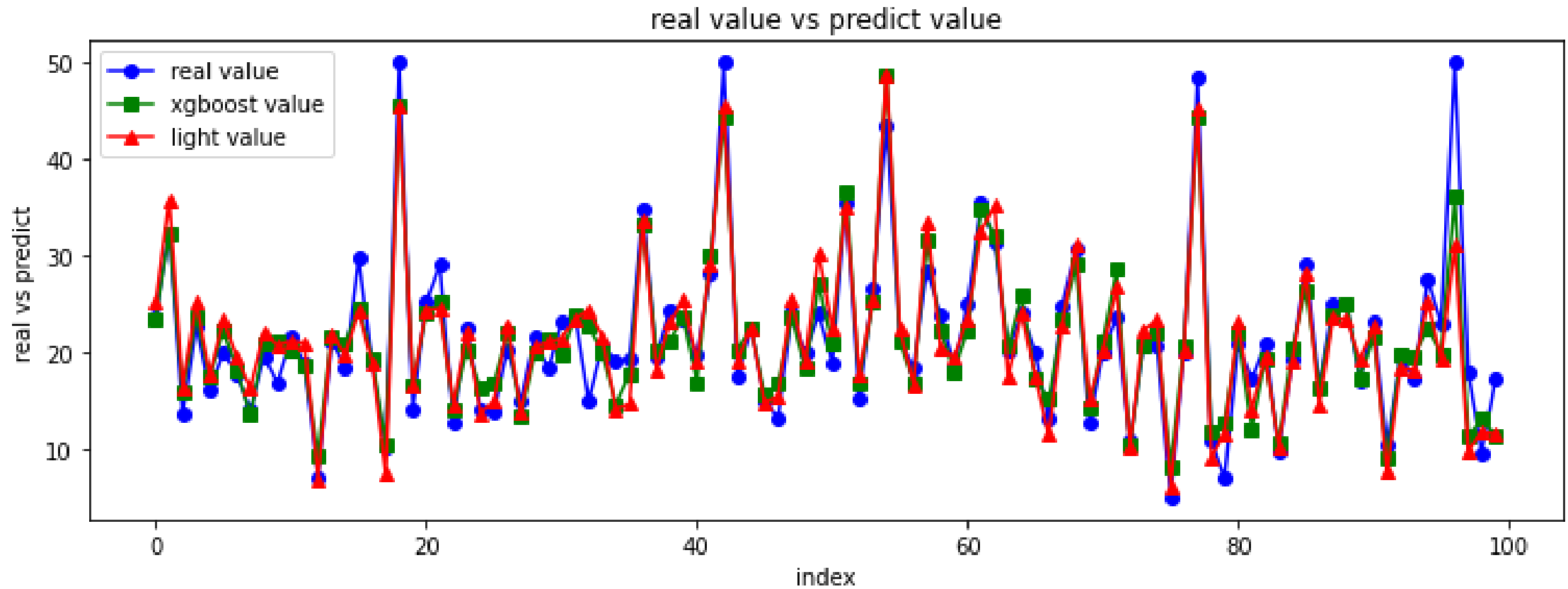
```
1 y_pred = model_lgb.predict(x_test)
2 y_true = y_test
```

```
1 mse = mean_squared_error(y_true, y_pred)
2 mse
```

11.346167217013914

1. Gradient Boosting

❖ LightGBM



1. Gradient Boosting

❖ Titanic

pclass : 1, 2, 3등석 정보를 각각 1, 2, 3으로 저장

survived 생존 여부. survived(생존), dead(사망)

name 이름

sex 성별. female(여성), male(남성)

age 나이

sibsp 함께 탑승한 형제 또는 배우자의 수

parch 함께 탑승한 부모 또는 자녀의 수

ticket 티켓 번호

fare 티켓 요금

cabin 선실 번호

embarked 탑승한 곳. C(Cherbourg), Q(Queenstown), S(Southampton)

1. Gradient Boosting

❖ Titanic

```
1 data=pd.read_csv("D:/big_data/titanic.csv")
2 #상위 데이터 보여주기
3 display(data.head())
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	S
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	NaN	S
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123	S
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	NaN	S

```
from google.colab import files
uploaded = files.upload()
import io
import pandas as pd

data = pd.read_csv(io.BytesIO(uploaded['titanic.csv']))

# data=pd.read_csv("D:/big_data/titanic.csv")
#상위 데이터 보여주기
display(data.head())
```

1. Gradient Boosting

❖ Titanic

```
1 print(data.shape)
2 data.columns.values
3
4 data.info()
```

(891, 12)

<class 'pandas.core.frame.DataFrame'>

RangeIndex: 891 entries, 0 to 890

Data columns (total 12 columns):

#	Column	Non-Null Count	Dtype
0	PassengerId	891 non-null	int64
1	Survived	891 non-null	int64
2	Pclass	891 non-null	int64
3	Name	891 non-null	object
4	Sex	891 non-null	object
5	Age	714 non-null	float64
6	SibSp	891 non-null	int64
7	Parch	891 non-null	int64
8	Ticket	891 non-null	object
9	Fare	891 non-null	float64
10	Cabin	204 non-null	object
11	Embarked	889 non-null	object

dtypes: float64(2), int64(5), object(5)

memory usage: 83.7+ KB

1. Gradient Boosting

❖ Titanic

```
1 y_data = data["Survived"]
2 data.drop(labels="Survived", axis=1, inplace=True)
3 display(data.head())
```

	PassengerId	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	S
1	2	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	NaN	S
3	4	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123	S
4	5	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	NaN	S

1. Gradient Boosting

❖ Titanic

```
1 drop_columns = ["PassengerId", "Name", "Age", "SibSp", "Ticket", "Cabin", "Parch", "Embarked"]
2 data.drop(labels=drop_columns, axis=1, inplace=True)
3 display(data.head())
```

	Pclass	Sex	Fare
0	3	male	7.2500
1	1	female	71.2833
2	3	female	7.9250
3	1	female	53.1000
4	3	male	8.0500

```
1 data = pd.get_dummies(data, columns=["Sex"])
2 data.fillna(value=0.0, inplace=True)
3 display(data.head())
```

	Pclass	Fare	Sex_female	Sex_male
0	3	7.2500	0	1
1	1	71.2833	1	0
2	3	7.9250	1	0
3	1	53.1000	1	0
4	3	8.0500	0	1

1. Gradient Boosting

❖ Titanic

```
1 state = 42
2 test_size = 0.30
3
4 X_train, X_val, y_train, y_val = train_test_split(data, y_data,
5     test_size=test_size, random_state=state)
```

1. Gradient Boosting

❖ Titanic

```
1 from sklearn.model_selection import GridSearchCV
2 from sklearn.ensemble import GradientBoostingClassifier
3
```

```
1
2 parameters = {
3     "learning_rate": [0.01, 0.025, 0.05, 0.075, 0.1, 0.15, 0.2],
4     "max_depth": [3, 5, 8],
5     "n_estimators": [10, 50, 100]
6 }
7
8
9
10 gb_cv = GridSearchCV(GradientBoostingClassifier(), parameters, cv=5, n_jobs=-1)
11 gb_cv.fit(X_train, y_train)
12
13 print('final params', gb_cv.best_params_) # 최적의 파라미터 값 출력
14 print('best score', gb_cv.best_score_)
15
16 cv_result_df=pd.DataFrame(gb_cv.cv_results_)
17 cv_result_df.sort_values(by=['rank_test_score'], inplace=True)
18 cv_result_df[['params', 'mean_test_score', 'rank_test_score']].head(10)
```

```
final params {'learning_rate': 0.2, 'max_depth': 5, 'n_estimators': 50}
best score 0.8218838709677421
```

	params	mean_test_score	rank_test_score
58	{'learning_rate': 0.2, 'max_depth': 5, 'n_esti...	0.821884	1
50	{'learning_rate': 0.15, 'max_depth': 5, 'n_est...	0.818658	2
59	{'learning_rate': 0.2, 'max_depth': 5, 'n_esti...	0.817032	3
48	{'learning_rate': 0.15, 'max_depth': 5, 'n_est...	0.817019	4
31	{'learning_rate': 0.075, 'max_depth': 5, 'n_es...	0.815458	5
44	{'learning_rate': 0.1, 'max_depth': 8, 'n_esti...	0.813871	6
26	{'learning_rate': 0.05, 'max_depth': 8, 'n_est...	0.813871	6
49	{'learning_rate': 0.15, 'max_depth': 5, 'n_est...	0.813858	8
35	{'learning_rate': 0.075, 'max_depth': 8, 'n_es...	0.813858	8
13	{'learning_rate': 0.025, 'max_depth': 5, 'n_es...	0.813832	10

1. Gradient Boosting

❖ Titanic

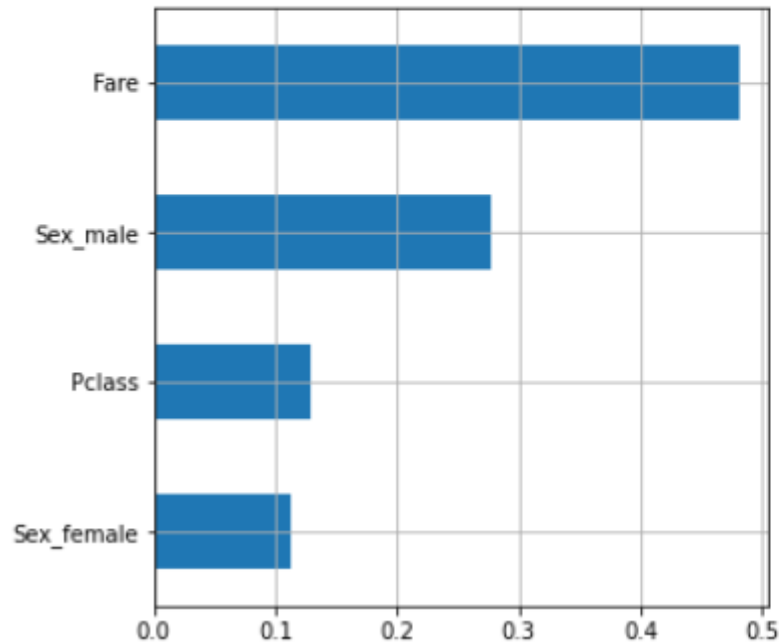
```
1 from sklearn.metrics import classification_report
2
3 gb_clf2 = GradientBoostingClassifier(n_estimators=50, learning_rate=0.2, max_depth=5, random_state=0)
4 gb_clf2.fit(X_train, y_train)
5 predictions = gb_clf2.predict(X_val)
6
7 print("Classification Report")
8 print(classification_report(y_val, predictions))
```

Classification Report				
	precision	recall	f1-score	support
0	0.81	0.90	0.85	157
1	0.83	0.69	0.75	111
accuracy			0.81	268
macro avg	0.82	0.80	0.80	268
weighted avg	0.81	0.81	0.81	268

1. Gradient Boosting

❖ Titanic

```
1  
2 feature_imp =pd.Series(gb_clf2.feature_importances_, index=data.columns).sort_values(ascending=True)  
3 feature_imp  
4 feature_imp.plot(kind='barh', grid=True, figsize=(5,5))  
5 plt.show()  
6
```



1. Gradient Boosting

❖ Titanic

```
1 from xgboost import plot_importance
2 from xgboost import XGBClassifier
3 import xgboost as xgb

1 xgg = XGBClassifier()
2
3
4 parameters = {
5     "learning_rate": [0.01, 0.025, 0.05, 0.075, 0.1, 0.15, 0.2],
6     "max_depth": [3, 5, 8],
7     "n_estimators": [10, 50, 100]
8 }
9
10
11
12 xgg_cv=GridSearchCV(xgg, param_grid=parameters, cv=5, scoring='accuracy', n_jobs=1)
13 xgg_cv.fit(X_train, y_train)
14
15 print('final params', xgg_cv.best_params_) # 최적의 파라미터 값 출력
16 print('best score', xgg_cv.best_score_)
17
18 cv_result_df=pd.DataFrame(xgg_cv.cv_results_)
19 cv_result_df.sort_values(by=['rank_test_score'], inplace=True)
20 cv_result_df[['params', 'mean_test_score', 'rank_test_score']].head(10)
```

```
final params {'learning_rate': 0.2, 'max_depth': 8, 'n_estimators': 100}
best score 0.8074322580645161
```

	params	mean_test_score	rank_test_score
62	{'learning_rate': 0.2, 'max_depth': 8, 'n_esti...	0.807432	1
40	{'learning_rate': 0.1, 'max_depth': 5, 'n_esti...	0.807368	2
57	{'learning_rate': 0.2, 'max_depth': 5, 'n_esti...	0.805768	3
23	{'learning_rate': 0.05, 'max_depth': 5, 'n_est...	0.805768	3
53	{'learning_rate': 0.15, 'max_depth': 8, 'n_est...	0.804232	5
14	{'learning_rate': 0.025, 'max_depth': 5, 'n_es...	0.804168	6
13	{'learning_rate': 0.025, 'max_depth': 5, 'n_es...	0.804168	6
22	{'learning_rate': 0.05, 'max_depth': 5, 'n_est...	0.804155	8
31	{'learning_rate': 0.075, 'max_depth': 5, 'n_es...	0.802555	9
46	{'learning_rate': 0.15, 'max_depth': 3, 'n_est...	0.800968	10

1. Gradient Boosting

❖ Titanic

```
1 xg_clf2 = XGBClassifier(n_estimators=100, learning_rate=0.2, max_depth=8, random_state=0)
2 X_train = pd.DataFrame(X_train, columns=data.columns)
3 X_test = pd.DataFrame(X_test, columns=data.columns)
4 # print(X_train)
5
6 xg_clf2.fit(X_train, y_train)
7 predictions = xg_clf2.predict(X_test)
8 print("Classification Report")
9 print(classification_report(y_test, predictions))
```

Classification Report				
	precision	recall	f1-score	support
0	0.81	0.90	0.85	157
1	0.83	0.69	0.75	111
accuracy			0.81	268
macro avg	0.82	0.80	0.80	268
weighted avg	0.81	0.81	0.81	268

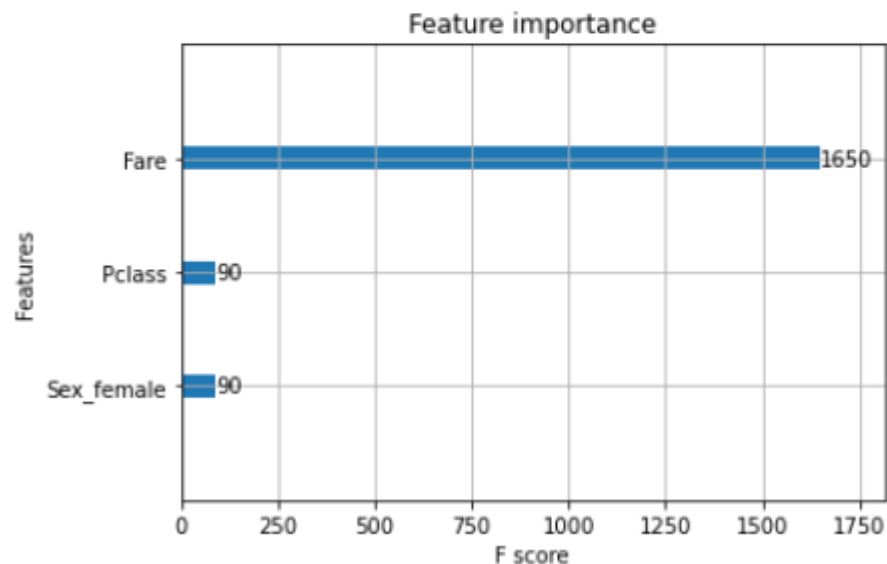
1. Gradient Boosting

❖ Titanic

```
1 fscore = xg_clf2.get_booster().get_fscore()
2 score_f= sorted(fscore.items(),key=(lambda x:x[1]), reverse=True)
3 print(score_f)
```

```
[('Fare', 1650), ('Sex_female', 90), ('Pclass', 90)]
```

```
1 import matplotlib.pyplot as plt
2 plot_importance(xg_clf2)
3 plt.show()
```



1. Gradient Boosting

❖ Titanic

```
1 from lightgbm import LGBMClassifier, plot_importance
2
3
4 lb = LGBMClassifier()
5
6 parameters = {
7     "learning_rate": [0.01, 0.025, 0.05, 0.075, 0.1, 0.15, 0.2],
8     "max_depth": [1, 2, 3, 4],
9     "n_estimators": [10, 50, 100]
10 }
11
12 lgb_cv=GridSearchCV(lb, param_grid=parameters, cv=5, scoring='accuracy', n_jobs=1)
13 lgb_cv.fit(X_train, y_train)
14
15 print('final params', lgb_cv.best_params_) # 최적의 파라미터 값 출력
16 print('best score', lgb_cv.best_score_)
17
18 cv_result_df=pd.DataFrame(lgb_cv.cv_results_)
19 cv_result_df.sort_values(by=['rank_test_score'], inplace=True)
20 cv_result_df[['params', 'mean_test_score', 'rank_test_score']].head(10)
```

final params {'learning_rate': 0.025, 'max_depth': 4, 'n_estimators': 50}
best score 0.8218580645161291

	params	mean_test_score	rank_test_score
81	{'learning_rate': 0.2, 'max_depth': 4, 'n_estimators': 100}	0.821858	1
22	{'learning_rate': 0.025, 'max_depth': 4, 'n_estimators': 50}	0.821858	1
69	{'learning_rate': 0.15, 'max_depth': 4, 'n_estimators': 100}	0.821858	1
11	{'learning_rate': 0.01, 'max_depth': 4, 'n_estimators': 100}	0.818658	4
57	{'learning_rate': 0.1, 'max_depth': 4, 'n_estimators': 100}	0.818658	4
46	{'learning_rate': 0.075, 'max_depth': 4, 'n_estimators': 100}	0.807394	6
34	{'learning_rate': 0.05, 'max_depth': 4, 'n_estimators': 100}	0.805858	7
23	{'learning_rate': 0.025, 'max_depth': 4, 'n_estimators': 100}	0.805858	7
70	{'learning_rate': 0.15, 'max_depth': 4, 'n_estimators': 100}	0.805794	9
64	{'learning_rate': 0.15, 'max_depth': 2, 'n_estimators': 100}	0.805781	10

1. Gradient Boosting

❖ Titanic

```
1 LG_clf2 = LGBMClassifier(n_estimators=50, learning_rate=0.025, max_depth=4, random_state=0)
2 X_train = pd.DataFrame(X_train, columns=data.columns)
3 X_test = pd.DataFrame(X_test, columns=data.columns)
4 # print(X_train)
5
6 LG_clf2.fit(X_train, y_train)
7 predictions = LG_clf2.predict(X_test)
8 print("Classification Report")
9 print(classification_report(y_test, predictions))
```

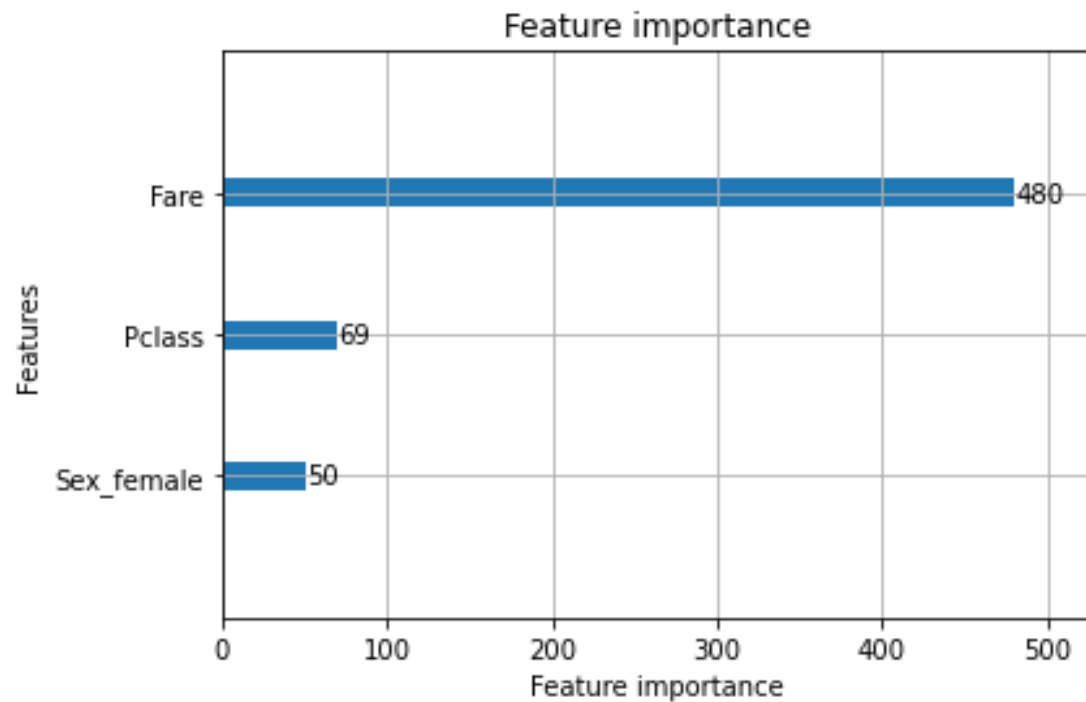
Classification Report

	precision	recall	f1-score	support
0	0.75	0.93	0.83	157
1	0.85	0.56	0.67	111
accuracy			0.78	268
macro avg	0.80	0.74	0.75	268
weighted avg	0.79	0.78	0.77	268

1. Gradient Boosting

❖ Titanic

```
1 import matplotlib.pyplot as plt
2 plot_importance(LG_clf2)
3 plt.show()
```



1. Gradient Boosting

❖ reference

- 모든 강의자료는 고려대 강필성 교수님 강의와 김성범 교수님 강의를 참고했음
- ratsgo's blog ,<https://ratsgo.github.io/>
- 안드레아스 뮐러, 세라 가이도 지음, 박해선 옮김, "파이썬 라이브러리를 활용한 머신러닝", 한빛미디어(2017)
- <https://injo.tistory.com/>
- <https://lsjsj92.tistory.com/>
- <https://scikit-learn.org/stable/index.html>

감사합니다