# 데이콘 유전체 정보 품종 예측

```python
1  # ars
2  ars=['SNP_02','SNP_03','SNP_04','SNP_09','SNP_11']
3  # bovine
4  bov=['SNP_05','SNP_06','SNP_08','SNP_15']
5  # hapmap
6  hap=['SNP_07','SNP_12','SNP_14']
7  # btb
8  btb=['SNP_10','SNP_13']
9
0  a.append('ars')
1  a.append('bov')
2  a.append('hap')
3  a.append('btb')
```

```python
# bovine

temp = ''

for i in bov:
    temp += train[i].str.replace(" ", "")

train['bov'] = temp


temp = ''

for i in bov:
    temp += test[i].str.replace(" ", "")

test['bov'] = temp
```

```python
1  # ars
2
3  temp = ''
4
5  for i in ars:
6      temp += train[i].str.replace(" ", "")
7
8  train['ars'] = temp
9
0
1  temp = ''
2
3  for i in ars:
4      temp += test[i].str.replace(" ", "")
5
6  test['ars'] = temp
```

```python
# hapmap

temp = ''

for i in hap:
    temp += train[i].str.replace(" ", "")

train['hap'] = temp


temp = ''

for i in hap:
    temp += test[i].str.replace(" ", "")

test['hap'] = temp
```
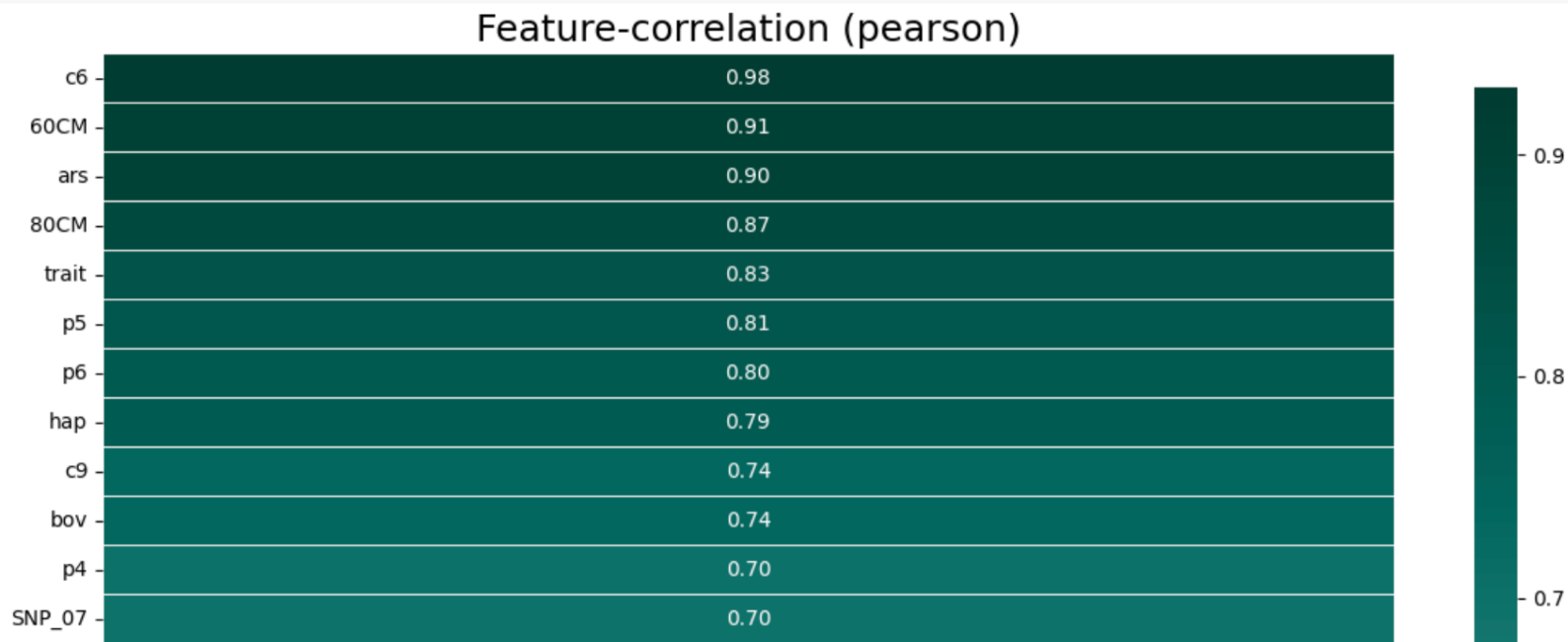
유전체 정보 피처 특징에 따라
피처 생성

```
1  ∨ klib.corr_plot(data = train.iloc[:,1:],
2  │   │   │   │          target = target,
3  │   │   │   │          figsize = (14, 12))
4
5  plt.show()
```

## Feature-correlation (pearson)

| Feature | Correlation |
|---------|-------------|
| c6 | 0.98 |
| 60CM | 0.91 |
| ars | 0.90 |
| 80CM | 0.87 |
| trait | 0.83 |
| p5 | 0.81 |
| p6 | 0.80 |
| hap | 0.79 |
| c9 | 0.74 |
| bov | 0.74 |
| p4 | 0.70 |
| SNP_07 | 0.70 |

상관관계 확인 및 정규화 작업

```
from sklearn.preprocessing import  StandardScaler


scaler = StandardScaler()

train.iloc[:,2:] = pd.DataFrame(scaler.fit_transform(train.iloc[:,2:]))
test.iloc[:,2:] = pd.DataFrame(scaler.fit_transform(test.iloc[:,2:]))
```

## ▾ XGBOOST

```
[95]  1   from xgboost import XGBClassifier

[96]  1   XGBClassifier

          xgboost.sklearn.XGBClassifier

[97]  1   select = SelectFromModel(estimator=XGBClassifier(), threshold='median')

[98]  1   select.fit(train.iloc[:,1:], target)
```

```
▸       SelectFromModel
▸ estimator: XGBClassifier
          ▸ XGBClassifier
```

```
[99]  1   # select 된 feature
      2   idx = select.get_support()
      3   train.iloc[:,1:].columns[idx]

          Index(['trait', 'SNP_01', 'SNP_07', 'SNP_10', 'SNP_13', 'SNP_14', 'SNP_15',
                 'ars', 'bov', 'hap', 'c6', '60CM', 'p4', 'p5'],
                dtype='object')

[100] 1   X_train = pd.DataFrame(select.transform(train.iloc[:,1:]))
      2
      3   y_train =  target
      4
      5   X_test = pd.DataFrame(select.transform(test.iloc[:,1:]))
```

```
[101] 1   import sklearn
      2   f1 = sklearn.metrics.f1_score

      1   from xgboost import XGBClassifier
      2
      3   m = XGBClassifier(
      4       booster='gbtree',
      5       max_depth=3,
      6       gamma=1,
      7       eta=0.4,
      8       reg_alpha=0.4,
      9       reg_lambda=0.8,
      10      min_child_weight=8,
      11
      12
      13      seed=777,
      14      n_estimators=100,
      15      colsample_bytree = 0.1,
      16      subsample = 0.4,
      17      objective='reg:linear',
      18      learning_rate=0.8,
      19      random_state=1
      20
      21  )
      22  m.fit(X_train, y_train)
      23  model_xg=pd.DataFrame(m.predict_proba(X_test))

[103] 1   model_xg= model_xg.reset_index()

[104] 1   model_xg['index'] = id_

[105] 1   model_xg.columns=['id','A','B','C']

[106] 1   model_xg['class']=model_xg.iloc[:,1:].idxmax(axis=1)
```

사용 모델 1. Xgboost

# LogitsicRegression

```
[107]  1   select = SelectFromModel(estimator=LogisticRegression(), threshold='median')
```

```
[108]  1   select.fit(train.iloc[:,1:], target)
```

```
          ▸         SelectFromModel
        ▸ estimator: LogisticRegression
                ▸ LogisticRegression
```

```
[109]  1   # select 된 feature
       2   idx = select.get_support()
       3   train.iloc[:,1:].columns[idx]
```

```
Index(['SNP_01', 'SNP_04', 'SNP_07', 'SNP_08'
       'SNP_15', 'ars', 'bov', 'c6', '60CM',
       dtype='object')
```

```
DataFrame: train
View
DataFrame with shape (262, 29)
```

```
       1   X_train = pd.DataFrame(select.transform(train.iloc[:,1:]))
       2
       3   y_train =  target
       4
       5   X_test = pd.DataFrame(select.transform(test.iloc[:,1:]))
```

```
[111]  1   from sklearn.linear_model import LogisticRegression
```

```
[112]  1   lrcv = LogisticRegression(random_state=0,
       2                            penalty='l2',
       3                            solver='saga',class_weight='balanced',C=0.1
       4                            )
```

```
       1   lrcv.fit(X_train, y_train)
```

```
              ▾          LogisticRegression
        LogisticRegression(C=0.1, class_weight='balanced', random_state=0,
                           solver='saga')
```

## ▸ XG+LR

```
[123]  1   a=[id_,((model_xg.iloc[:,1:4]*0.5 + model_Logistic.iloc[:,1:4]*0.5))]
```

```
[124]  1   ensemble_proba=pd.DataFrame(pd.concat(a,axis=1))
```

```
[125]  1   ensemble=pd.DataFrame(pd.concat(a,axis=1).iloc[:,1:].idxmax(axis=1))
```

```
[126]  1   ensemble.reset_index(inplace=True)
```

```
[127]  1   ensemble.columns=['id','class']
```

```
[128]  1   ensemble.id = id_
```

```
[ ]    1   ensemble.to_csv("XG+LR.csv",index=False)
```

사용 모델 2. Logistic Regression

## ▾ XG+LR

```
[123]  1  a=[id_,((model_xg.iloc[:,1:4]*0.5 + model_Logistic.iloc[:,1:4]*0.5))]
```

```
✓ [124]  1  ensemble_proba=pd.DataFrame(pd.concat(a,axis=1))
```

```
✓ [125]  1  ensemble=pd.DataFrame(pd.concat(a,axis=1).iloc[:,1:].idxmax(axis=1))
```

```
✓ [126]  1  ensemble.reset_index(inplace=True)
```

```
✓ [127]  1  ensemble.columns=['id','class']
```

```
✓ [128]  1  ensemble.id = id_
```

```
[ ]  1  ensemble.to_csv("XG+LR.csv",index=False)
```

두 모델의 결과 앙상블

# 유전체 정보 품종 분류 AI 경진대회

알고리즘 | 유전체 | 분류 | Macro F1 Score

₩ 상금 : 300 만원

🕐 2022.12.12 ~ 2023.01.16 09:59　+ Google Calendar

👥 1,316명　📅 마감

🔗

**참여중**

---

대회안내　데이터　코드 공유　토크　**리더보드**　제출

---

PUBLIC　**PRIVATE**　AWARDS　RANKING CHART

순위기준

● WINNER　● 1%　● 4%　● 10%　전체 랭킹 ›

| # | 팀 | 팀 멤버 | 최종점수 | 제출수 | 등록일 |
|---|---|---|---|---|---|
| 140 | 헬응애 | 🟦 헬응 | 0.9673 | 54 | 9달 전 |

Private date 데이터
최종 140위
(상위 20%) 달성