# Introduction

The goal of this notebook is to provide some technical details (key steps) for our submitted manuscript:

> *Haolin Wang, et al. "Integrating Co-clustering and Interpretable Machine Learning for the Prediction of Intravenous Immunoglobulin Resistance in Kawasaki Disease", 2020.*

To enable clinically applicable prediction addressing the incompleteness of clinical data and the lack of interpretability of machine learning models, a multi-stage method is developed by integrating data missing pattern mining and intelligible models. First, co-clustering is adopted to characterize the block-wise data missing patterns by simultaneously grouping the clinical features and patients to enable (a) group-based feature selection and missing data imputation and (b) patient subgroup-specific predictive models considering the availability of data. Second, feature selection is performed using the group Lasso to uncover group-specific risk factors. Third, the Explainable Boosting Machine, which is an interpretable learning method based on generalized additive models, is applied for the prediction of each patient subgroup.
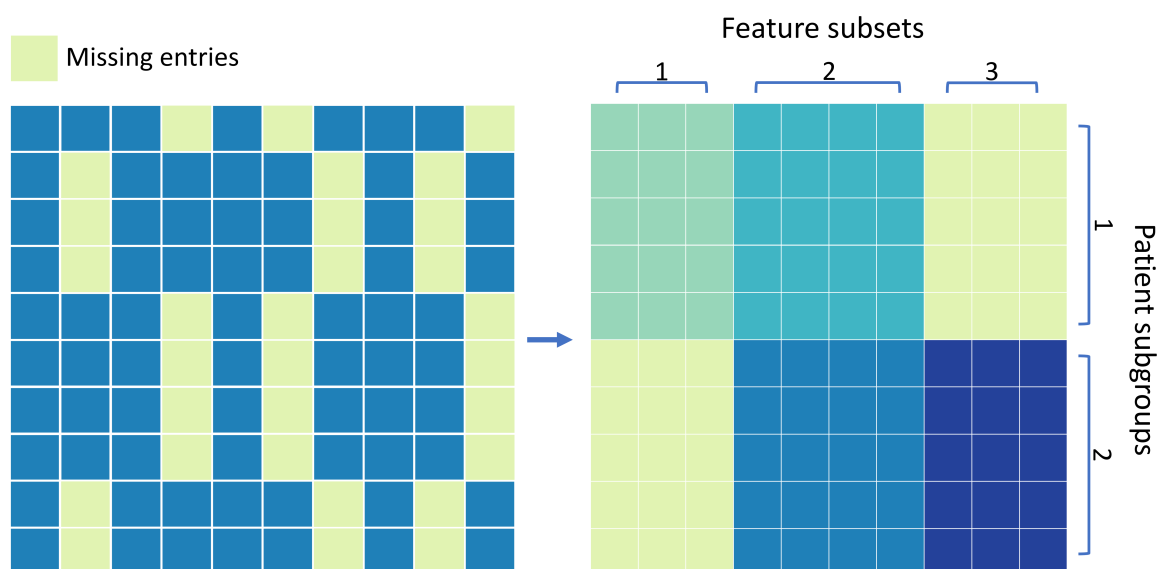


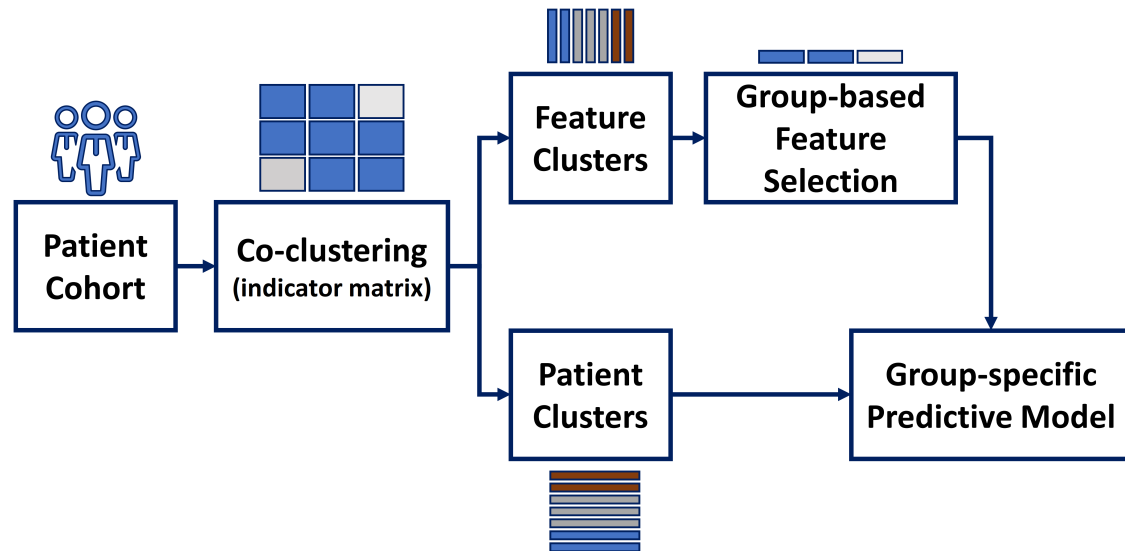Fig 1. The block-wise missing patterns characterized by co-clustering.

Fig 2. The proposed multiple classifier system with static classifier selection based on co-clustering.

## Tools and Implementation

### Packages

- Popular Python machine learning libraries such as scikit-learn and pandas
- imbalanced-learn: https://github.com/scikit-learn-contrib/imbalanced-learn
- InterpretML: https://github.com/interpretml/interpret
- Coclust: https://pypi.org/project/coclust/
- Group Lasso: https://github.com/AnchorBlues/GroupLasso/blob/master/grouplasso/model.py

### Preprocessing

```
1  import pandas as pd
2
3  dat = pd.read_csv('dataset.csv')
4  df = pd.get_dummies(dat, columns=['categories'])
```

### Over-sampling for imbalanced dataset

```
1  from imblearn.over_sampling import SMOTE
2
3  balanced_feature_set, balanced_label = SMOTE(sampling_strategy=<>,
   random_state=0).fit_resample(feature, label)
```

### Missing data imputation to address data missing at random

```
1  from sklearn.experimental import enable_iterative_imputer
2  from sklearn.impute import IterativeImputer
3
4  imp = IterativeImputer(max_iter=5, random_state=0, tol=0.005)
5  feature_set = imp.fit_transform(feature_set)
```

### Dataset splitting for cross-validation

```
1  from sklearn.model_selection import StratifiedKFold
2
3  skf = StratifiedKFold(n_splits=5)
4  skf.get_n_splits(X, y)
5  for train_index, test_index in skf.split(X, y):
6      X_train, X_test = X[train_index], X[test_index]
7      y_train, y_test = y[train_index], y[test_index]
```

## Baseline methods

```
1   #Lasso
2   from sklearn.linear_model import Lasso
3   #Logistic
4   from sklearn.linear_model import LogisticRegression
5   #Ridge
6   from sklearn.linear_model import Ridge
7   #KNN
8   from sklearn import neighbors
9   model = neighbors.KNeighborsClassifier()
10  #Naive Bayes
11  from sklearn.naive_bayes import GaussianNB
12  #MLP
13  from sklearn.neural_network import MLPClassifier
14  #Random forest
15  from sklearn.ensemble import RandomForestClassifier
16  #lightGBM
17  import lightgbm
18  #XGBoost
19  import xgboost as xgb
20  #EBM
21  from interpret.glassbox import ExplainableBoostingClassifier
22
23  #Those models are trained and tested using the same splitted dataset for
    cross-validation. The usage of lightGBM is slightly different.
24  train_data = lgb.Dataset(train_set, label=train_label)
25  param = {'metrics':'auc', 'objective': 'binary'}
26  model = lgb.train(param, train_data)
27  prob = model.predict(test_set)
```

## Tuning the hyper-parameters

```
1  #alpha for Lasso
2  from sklearn.linear_model import Lasso
3  from sklearn.model_selection import GridSearchCV
4
5  params = {"alpha": numpy.logspace(-3, 1, 5)}
6  model_cv = GridSearchCV(Lasso(), params, cv=5)
7  model = model_cv.fit(train_set, train_label)
8  print("tuned hpyerparameters:", model.best_params_)
9  prob = model_cv.predict(test_set)
```

## The AUC score of the five-fold cross-validation of the baseline methods

|          | 1     | 2     | 3     | 4     | 5     | MEAN  | STD   |
|----------|-------|-------|-------|-------|-------|-------|-------|
| Lasso    | 0.800 | 0.847 | 0.880 | 0.847 | 0.864 | 0.848 | 0.027 |
| Logistic | 0.802 | 0.847 | 0.888 | 0.856 | 0.874 | 0.853 | 0.029 |
| Ridge    | 0.811 | 0.865 | 0.897 | 0.849 | 0.876 | 0.860 | 0.029 |
| KNN      | 0.837 | 0.801 | 0.813 | 0.824 | 0.827 | 0.820 | 0.012 |
| MNB      | 0.855 | 0.866 | 0.867 | 0.826 | 0.830 | 0.849 | 0.017 |
| MLP      | 0.809 | 0.870 | 0.858 | 0.818 | 0.822 | 0.836 | 0.024 |
| RF       | 0.865 | 0.849 | 0.866 | 0.848 | 0.846 | 0.855 | 0.009 |
| EBM      | 0.882 | 0.875 | 0.878 | 0.869 | 0.888 | 0.878 | 0.006 |
| GBM      | 0.880 | 0.871 | 0.885 | 0.875 | 0.909 | 0.884 | 0.013 |
| XGB      | 0.884 | 0.876 | 0.886 | 0.866 | 0.900 | 0.882 | 0.011 |

## Co-clustering and Classification

### Indicator matrix for incomplete dataset

```
1   # generate indicator matrix using the original dataset without missing data
    imputation
2   [rows, cols] = train_set_original.shape
3   train_ind = numpy.zeros((rows, cols))
4   for r in range(rows):
5       for c in range(cols):
6           if not numpy.isnan(train_set_original[r, c]):
7               train_ind[r, c] = 1
8   train_rows_num = rows
9
10  ...
11
12  # merge training set and test set splitted for cross-validation
13  full_ind = numpy.row_stack((train_ind, test_ind))
```

### Co-clustering of the indicator matrix

```
1   from coclust.coclustering import CoclustInfo
2
3   for row_cluster in range(2, 10):
4       for column_cluster in range(2,30):
5           ...
6           model =
    CoclustInfo(n_row_clusters=row_cluster,n_col_clusters=column_cluster,
    random_state=42)
7           model.fit(full_ind)
8           row_labels = model.row_labels_
9           print(row_labels)
10          col_labels = model.column_labels_
11          print(col_labels)
12          ...
```

**Visualization**

```
1   row_indices = numpy.argsort(model.row_labels_)
2   col_indices = numpy.argsort(model.column_labels_)
3   X_reorg = full_ind[row_indices, :]
4   X_reorg = X_reorg[:, col_indices]
5   cmap = sns.color_palette("YlGnBu", 41)
6   fig = sns.heatmap(X_reorg, cmap=cmap, xticklabels=False, yticklabels=False)
7   plt.savefig('co-cluster.tif', dpi=600, format='tif')
```
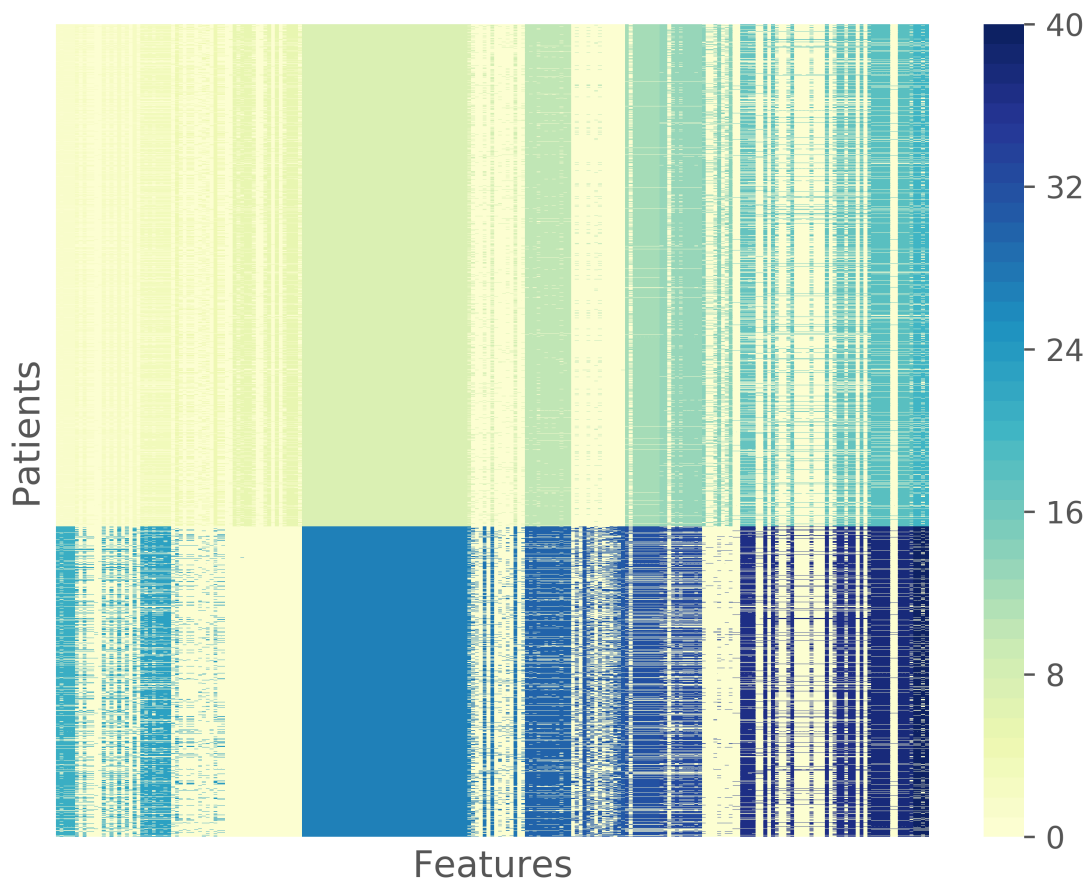


Fig 3. Co-clustering of the real-world clinical dataset

**Re-organize samples for each row clusters to train multiple classifiers**

```
1   for cluster_index in range(0, row_cluster):
2       co_train_set = []
3       co_train_label = []
4       for k in range(0, train_rows_num):
5           if row_labels[k] == cluster_index:
6               co_train_label.append(train_label[k])
7               co_train_set.append(train_set[k,:])
8       co_train_set = numpy.array(co_train_set)
9       co_train_label = numpy.array(co_train_label)
10      print(co_train_set.shape)
11      print(co_train_label.shape)
```

## Feature selection with group feature structure

Group Lasso derives feature coefficients from certain groups to be small or exact zero.

```
1   # train group lasso
2   from grouplasso import GroupLassoClassifier
3
4   model =
    GroupLassoClassifier(group_ids=numpy.array(col_labels),alpha=alpha_test,
    eta=0.001, tol=0.001, max_iter=3000, random_state=42, verbose_interval=300)
5   model.fit(co_train_set, co_train_label)
6   print(model.coef_)
7
8   # feature selection
9   selected_cols = []
10  for k in range(0, len(col_labels)):
11      if abs(model.coef_[k]) > 0:
12          selected_cols.append(k)
```

## Train group-specific prediction model

```
1   decision_model = ExplainableBoostingClassifier().fit(co_train_set[:,
    selected_cols], co_train_label)
2   predict_test += list(decision_model.predict_proba(co_test_set[:,
    selected_cols])[:,1])
3   predict_test_label += list(co_test_label)
```

## Evaluation metrics

```
1   from sklearn.metrics import precision_recall_fscore_support
2
3   #for 5-fold cross-validation
4   for index in range(1, 6):
5       best_f1 = 0
6       keep_score = []
7       test_label = np.load('xxx.npy')
8       pred = np.load('xxx.npy')
9       fpr, tpr, thresholds = metrics.roc_curve(test_label, prob)
10      roc_auc = metrics.auc(fpr, tpr)
11      print(roc_auc)
12
13      for thres in pred:
14          res_bin = []
```

```
15          for p in prob:
16              if p >= thres:
17                  res_bin.append(1)
18              else:
19                  res_bin.append(0)
20
21          score = precision_recall_fscore_support(test_label, res_bin,
    average='binary')
22          print(score)
23          if score[2] > best_f1:
24              keep_score = score
25              best_f1 = score[2]
26      best_score.append(keep_score)
27      average_pr.append((keep_score[0] + keep_score[1])/2)
28
29  print(best_score)
30  best_score_array = np.array(best_score)
31  average_pr = np.array(average_pr)
32
33  print( str(np.around(best_score_array[:,0].mean(), decimals=3)) + '+' +
    str(np.around(best_score_array[:,0].std(), decimals=3)))
34  print( str(np.around(best_score_array[:,1].mean(), decimals=3)) + '+' +
    str(np.around(best_score_array[:,1].std(), decimals=3)))
35  print( str(np.around(best_score_array[:,2].mean(), decimals=3)) + '+' +
    str(np.around(best_score_array[:,2].std(), decimals=3)))
36  print( str(np.around(average_pr.mean(), decimals=3)) + '+' +
    str(np.around(average_pr.std(), decimals=3)))
```

**The AUC score of the five-fold cross-validation of the models enhanced by the proposed framework.**

|          | 1       | 2       | 3       | 4       | 5       | MEAN    | STD   |
|----------|---------|---------|---------|---------|---------|---------|-------|
| Lasso    | 0.804   | 0.891   | 0.897   | 0.859   | 0.869   | 0.864   | 0.033 |
| Logistic | 0.835   | 0.887   | 0.875   | 0.878   | 0.874   | 0.87    | 0.018 |
| Ridge    | 0.873   | 0.894   | 0.896   | 0.884   | 0.864   | 0.882   | 0.012 |
| KNN      | 0.841   | 0.827   | 0.82    | 0.832   | 0.829   | 0.83    | 0.007 |
| MNB      | 0.84    | 0.883   | 0.858   | 0.861   | 0.872   | 0.863   | 0.014 |
| MLP      | 0.791   | 0.881   | 0.855   | 0.834   | 0.848   | 0.842   | 0.03  |
| RF       | 0.869   | 0.884   | 0.859   | 0.873   | 0.875   | 0.872   | 0.008 |
| GBM      | **0.888** | 0.892 | 0.925   | 0.908   | 0.9     | 0.903   | 0.013 |
| XGB      | **0.888** | 0.891 | 0.926   | 0.909   | 0.899   | 0.903   | 0.014 |
| EBM      | 0.88    | **0.919** | **0.946** | **0.916** | **0.923** | **0.917** | 0.021 |

## ROC curve (mean+std)

```
1  import matplotlib.pyplot as plt
2  import numpy
```

```
 3   from scipy import interp
 4
 5   plt.style.use('seaborn-white')
 6
 7   model_tprs = []
 8   model_aucs = []
 9   model_mean_fpr = []
10   model_mean_fpr = numpy.linspace(0, 1, 100)
11
12   model_fpr = numpy.load("../proposed/temp/model_" + str(group) + "_fpr.npy")
13   model_tpr = numpy.load("../proposed/temp/model_" + str(group) + "_tpr.npy")
14   model_auc = numpy.load("../proposed/temp/model_" + str(group) + "_auc.npy")
15   model_tprs.append(interp(model_mean_fpr, model_fpr, model_tpr))
16   model_tprs[-1][0] = 0.0
17   model_aucs.append(model_auc)
18
19   model_mean_tpr = numpy.mean(model_tprs, axis=0)
20   model_mean_tpr[-1] = 1.0
21   model_mean_auc = numpy.mean(model_aucs)
22   model_std_auc = numpy.std(model_aucs)
23   plt.plot(model_mean_fpr, model_mean_tpr, color='b',
24            label=r'EBM (AUC = %0.3f $\pm$ %0.3f)' % (model_mean_auc,
     model_std_auc),
25            lw=2, alpha=.8)
26   model_std_tpr = numpy.std(model_tprs, axis=0)
27   model_tprs_upper = numpy.minimum(model_mean_tpr + model_std_tpr, 1)
28   model_tprs_lower = numpy.maximum(model_mean_tpr - model_std_tpr, 0)
29   # label=r'$\pm$ 1 std. dev.'
30   plt.fill_between(model_mean_fpr, model_tprs_lower, model_tprs_upper,
     color='b', alpha=.2)
```
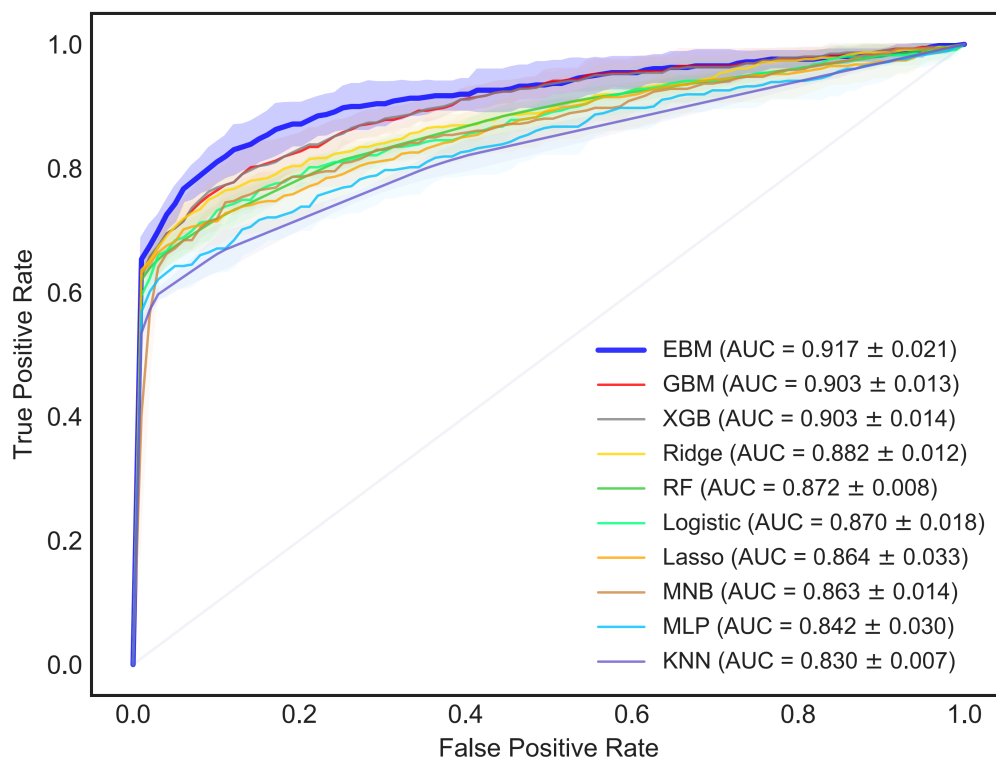


Fig 4. ROC curves

## Interpretability

```
 1  model = ExplainableBoostingClassifier()
 2  model.fit(train_set, train_label)
 3
 4  ebm_global = model.explain_global()
 5  # show(ebm_global)
 6  # export model parameters
 7  pd.DataFrame(ebm_global._internal_obj['overall']).to_csv('xxx.csv',
    index=False, header=False)
 8
 9  ebm_local = model.explain_local(train_set, train_label)
10  # export model parameters
11  pd.DataFrame(ebm_local._internal_obj['specific'][<sample_id>]
    ['scores']).to_csv('xxx.csv', index=False, header=False)
```