

XGBoost Basics

Learn about the basics of using XGBoost.

Chapter Goals:

- Learn about the XGBoost data matrix
- Train a `Booster` object in XGBoost

A. Basic data structures

The basic data structure for XGBoost is the `DMatrix`, which represents a data matrix. The `DMatrix` can be constructed from NumPy arrays.

The code below creates `DMatrix` objects with and without labels.

```
1 data = np.array([
2     [1.2, 3.3, 1.4],
3     [5.1, 2.2, 6.6]])
4
5 import xgboost as xgb
6 dmat1 = xgb.DMatrix(data)
7
8 labels = np.array([0, 1])
9 dmat2 = xgb.DMatrix(data, label=labels)
```

RUN

SAVE

RESET



The `DMatrix` object can be used for training and using a `Booster` object, which represents the gradient boosted decision tree. The `train` function in XGBoost lets us train a `Booster` with a specified set of parameters.

The code below trains a `Booster` object using a predefined labeled dataset.

```
1 # predefined data and labels
2 print('Data shape: {}'.format(data.shape))
3 print('Labels shape: {}'.format(labels.shape))
4 dtrain = xgb.DMatrix(data, label=labels)
5
6 # training parameters
7 params = {
8     'max_depth': 0,
9     'objective': 'binary:logistic'
10 }
11 print('Start training')
12 bst = xgb.train(params, dtrain) # booster
13 print('Finish training')
```

RUN

SAVE

RESET



Close

Output

2.728s

```
Data shape: (569, 30)
Labels shape: (569,)
Start training
[16:46:47] /workspace/src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[16:46:47] /workspace/src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[16:46:47] /workspace/src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[16:46:47] /workspace/src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[16:46:47] /workspace/src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[16:46:47] /workspace/src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
```

A list of the possible parameters and their values can be found [here](#). In the example above, we set the `'max_depth'` parameter to `0` (which means no limit on the tree depths, equivalent to `None` in scikit-learn). We also set the `'objective'` parameter (the objective function) to binary classification via logistic regression. For the remaining available parameters, we used their default settings (so we didn't include them in `params`).

B. Using a **Booster**

After training a **Booster**, we can evaluate it and use it to make predictions.

```
1 # predefined evaluation data and labels
2 print('Data shape: {}'.format(eval_data.shape))
3 print('Labels shape: {}'.format(eval_labels.shape))
4 deval = xgb.DMatrix(eval_data, label=eval_labels)
5
6 # Trained bst from previous code
7 print(bst.eval(deval)) # evaluation
8
9 # new_data contains 2 new data observations
10 dpred = xgb.DMatrix(new_data)
11 # predictions represents probabilities
12 predictions = bst.predict(dpred)
13 print('{}\n'.format(predictions))
```

RUN

SAVE

RESET



Close

Output

3.596s

```
Data shape: (119, 30)
Labels shape: (119,)
[0] eval-error:0.226891
[0.6236573 0.6236573]
```

The evaluation metric used for binary classification (**eval-error**) represents the classification error, which is the default **'eval_metric'** parameter for binary classification **Booster** models.

Note that the model's predictions (from the **predict** function) are probabilities, rather than class labels. The actual label classifications are just the rounded probabilities. In the example above, the **Booster** predicts classes of 0 and 1, respectively.