

需要给输入数据指定相应的格式。

在TVM中，特征通过C++相关代码获取并打包，在Python进行解包并用于xgboost模型的更新。

```
1 features, normalized_throughputs, task_ids =
2     get_per_store_features_from_measure_pairs(
3     self.inputs, self.results, skip_first_n_feature_extraction=n_cache
```

即输入数据包含“特征”，“标准化吞吐量”和“任务id”。

然后将这几种数据打包为xgb.DMatrix格式，并按照task id按顺序：

```
1 dtrain = pack_sum_xgbmatrix(
2     features, normalized_throughputs, task_ids, normalized_throughput
3 )
```

在pack_sum_xgbmatrix () 中，通过一系列子任务排序等操作，最终将(feature, label)定义xgb.DMatrix格式的数据。

```
1 ret = xgb.DMatrix(np.array(x_flatten), y_flatten)
```

至此完成xgboost模型训练数据的格式转换。

参考TVM的做法，我们没有子任务的划分，即对于整个程序，我们提取出与运行时间有关的特征以及程序相应的运行时间即可。

如何创造训练数据

由于没有TVM那样的代码生成和进化搜索过程，我们肯定会面临样本不足的问题。

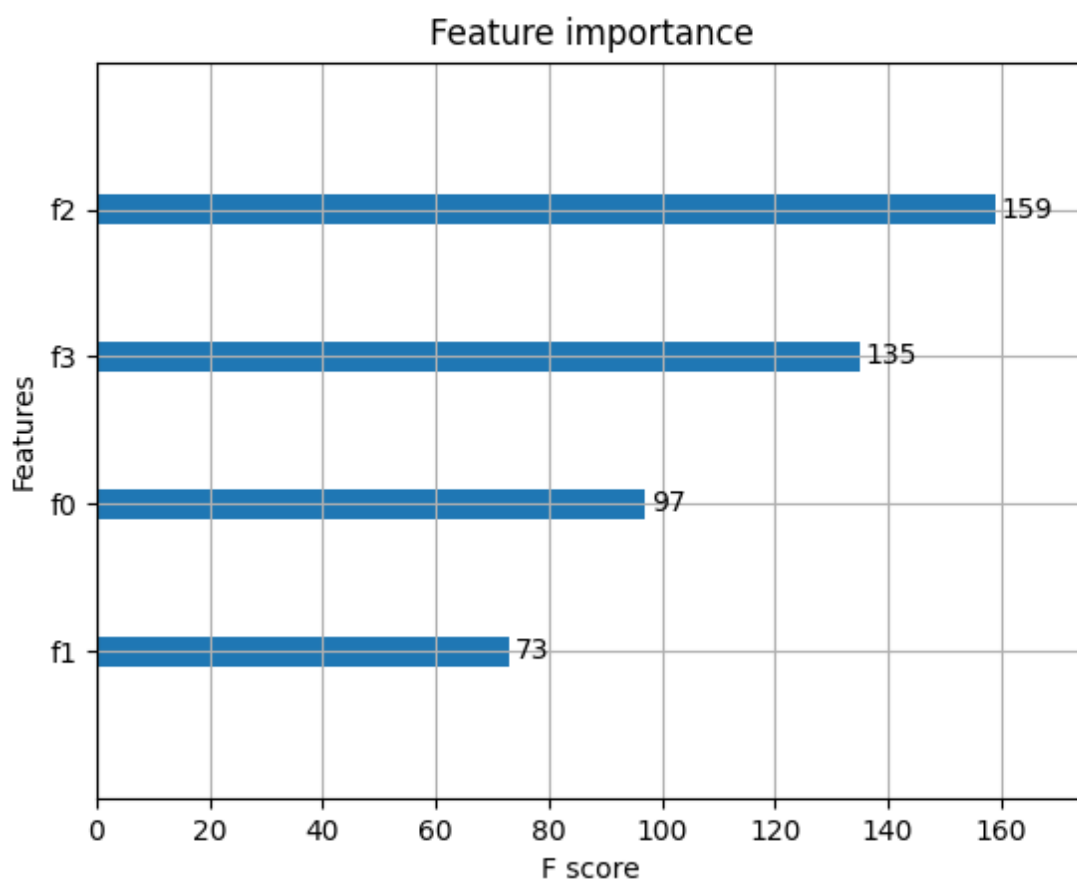
但是如何生成最佳代码不是cost model的考虑，它所要做的就是对代码的评估，理论上只要我们微调代码，跑出runtime，得到足够多的程序特征-运行时间对，我们就可以拟合出一个可以根据特征预测运行时间的回归树模型。

如果获取到的不是groundtruth，而是根据规则计算出来的“runtime”，那么模型最终拟合出来的也只能是规则本身。

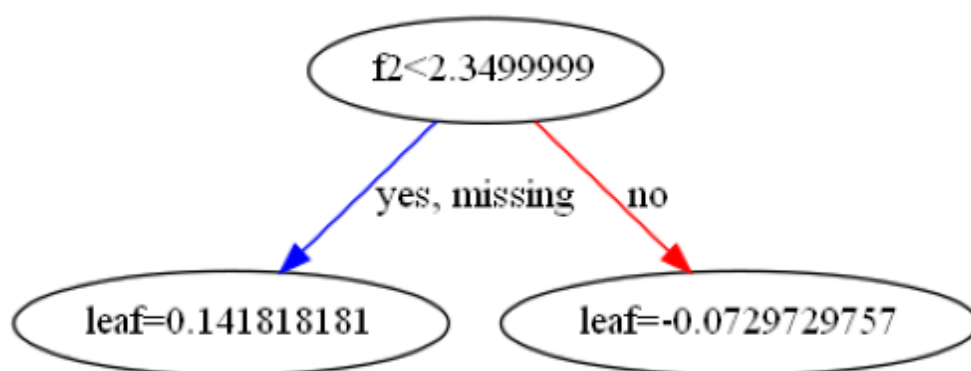
如何调用XGBoost模型

目前还没有特征提取的接口，那么特征数据就是裸露的文件，用C++或Python都是可以读取的。

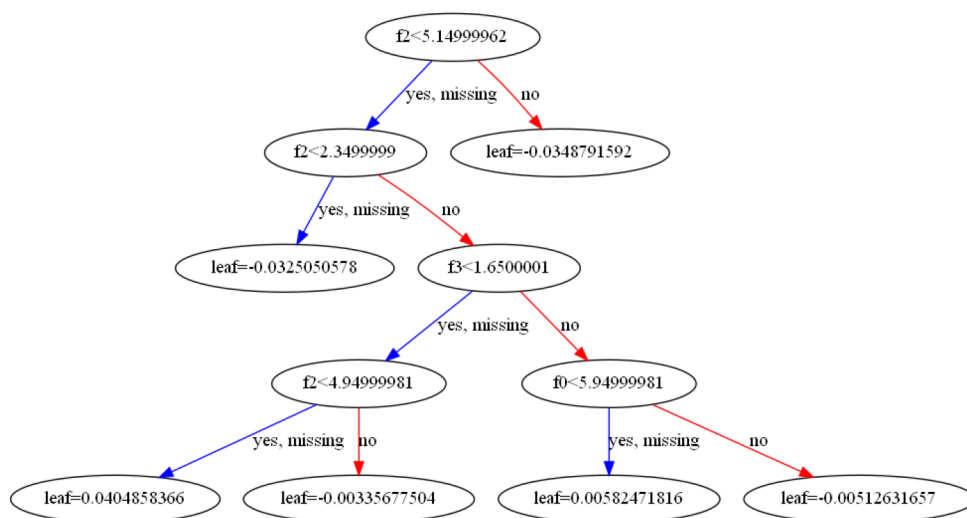
下面是XGBoost完成训练后，对不同的特征的重要性排序（值越大说明特征越重要，所在树的节点也越上层）：



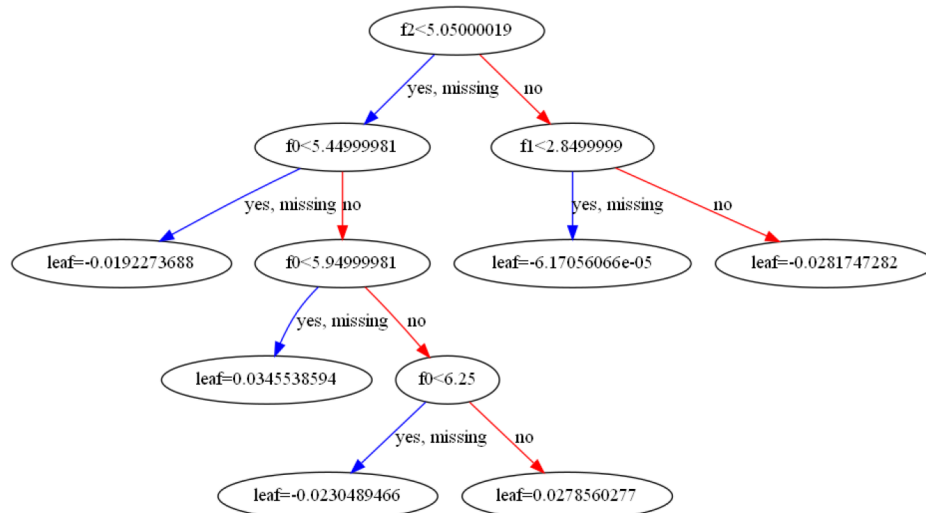
训练结果中的决策树：
num_trees=0第一棵树



num_trees=100



num_trees=130, 有stop_earlying的参数



下图是xgb实战相关

