

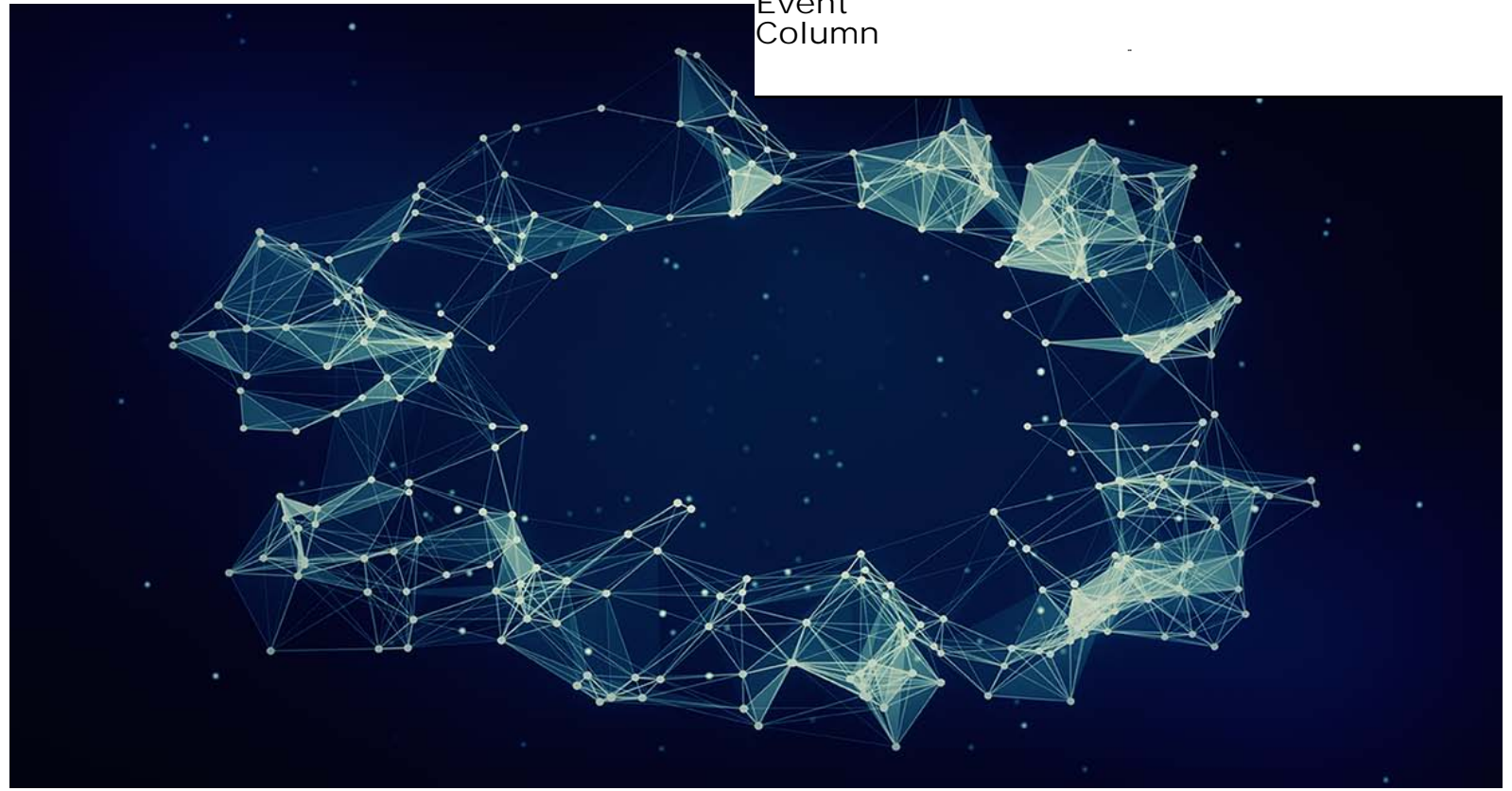
手を動かして GBDT を理解

t2sy

ツイート

0

NHN TECHORUS
Tech Blog
AWS
Data Science
Tech
Event
Column



- Topics
- GBDT
 - おわりに

こんにちは。データサイエンスチームの t2sy です。
この記事では、多くの機械学習タスクで使われている GBDT (Gradient Boosting Decision Tree) を手を動かして実装・実験することでアルゴリズムを理解することを目指します。
環境は Python 3.7.2、NumPy 1.15.4 です。

GBDT

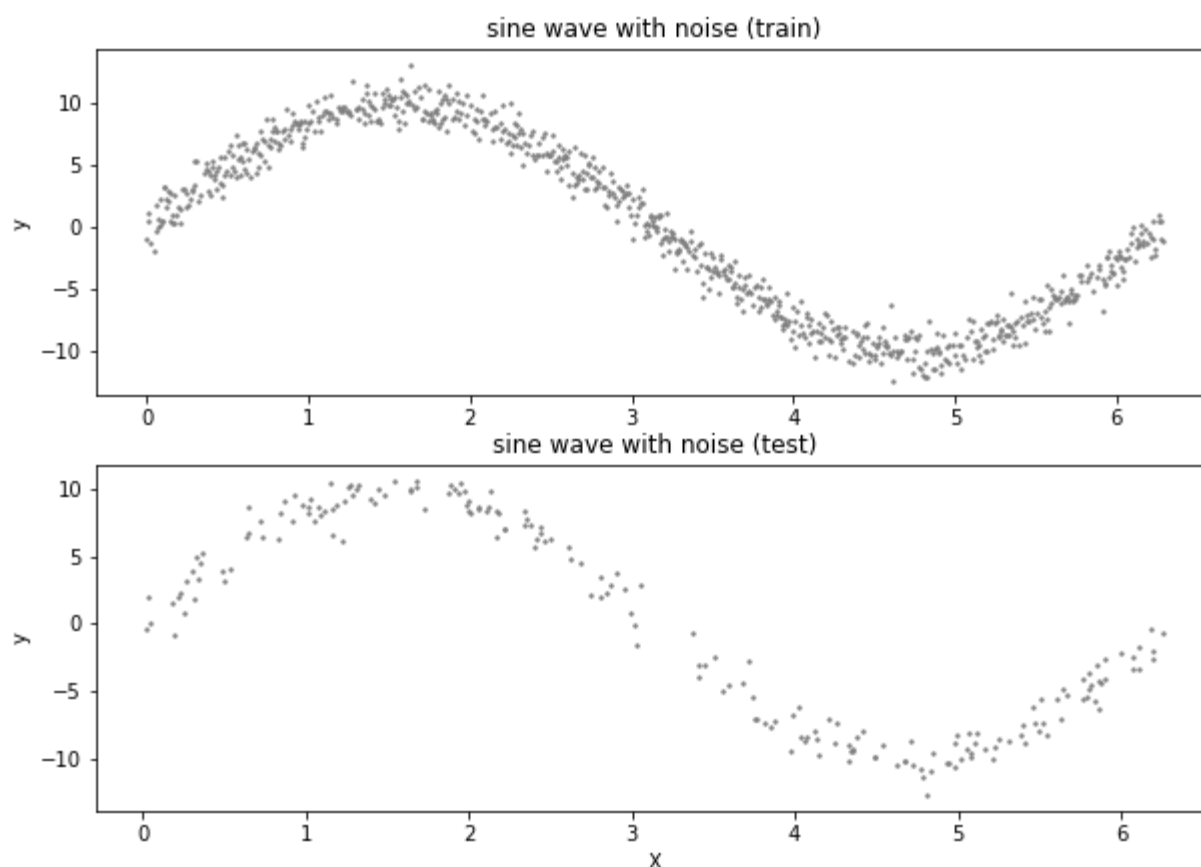
ブースティング (Boosting) は弱学習器 (weak learner) の学習を逐次的に行い強学習器に導くアルゴリズムを指します。GBDT は損失関数の勾配を用いたブースティングで弱学習器に決定木を用います。代表的な GBDT の実装には以下があります。

- [XGBoost](#)
- [LightGBM](#)
- [CatBoost](#)
- [scikit-learn](#)

データ

最初の実験で使うデータを準備します。振幅10の正弦波に標準正規分布 $N(0,1)$ を乗せた1000行1列のデータを生成し、訓練:テスト=8:2で分割します。

```
1 import numpy as np
2 from sklearn.model_selection import train_test_split
3
4 def generate_data(N=1000, test_size=0.2, seed=123):
5     np.random.seed(seed)
6     X = np.linspace(0, 2 * np.pi, N)
7     X = X.reshape(-1, 1)
8     y = 10 * np.sin(X[:, 0]) + np.random.standard_normal(N)
9
10    return train_test_split(X, y, test_size=test_size, random_state=seed)
```



生成したデータに対して回帰木と GBDT で曲線フィッティングを行い比較します。

決定木

GBDT を構成する弱学習器である決定木あるいは回帰木は、情報利得が最大となる特徴でデータを再帰的に分割するアルゴリズムです。通常は二分決定木となります。分割条件は分類問題の場合はエントロピー、ジニ不純度、分類誤差など、回帰問題の場合は、MSE (mean squared error)、LAD (least absolute deviation) などがあります。

まず、Tree クラスを書きます。

```
1 class Tree(object):
2     def __init__(self, pre_pruning=False, max_depth=6):
3         self.feature = None
```

```

4     self.label = None
5     self.n_samples = None
6     self.gain = None
7     self.left = None
8     self.right = None
9     self.threshold = None
10    self.pre_pruning = pre_pruning
11    self.max_depth = max_depth
12    self.depth = 0

13    def build(self, features, target, criterion='gini'):
14        self.n_samples = features.shape[0]

15        if len(np.unique(target)) == 1:
16            self.label = target[0]
17            return

18        best_gain = 0.0
19        best_feature = None
20        best_threshold = None

21        if criterion in {'gini', 'entropy', 'error'}:
22            self.label = max(target, key=lambda c: len(target[target==c]))
23        else:
24            self.label = np.mean(target)

25        impurity_node = self._calc_impurity(criterion, target)

26        for col in range(features.shape[1]):
27            feature_level = np.unique(features[:,col])
28            thresholds = (feature_level[:-1] + feature_level[1:]) / 2.0

29            for threshold in thresholds:
30                target_l = target[features[:,col] <= threshold]
31                impurity_l = self._calc_impurity(criterion, target_l)
32                n_l = target_l.shape[0] / self.n_samples

33                target_r = target[features[:,col] > threshold]
34                impurity_r = self._calc_impurity(criterion, target_r)
35                n_r = target_r.shape[0] / self.n_samples

36                ig = impurity_node - (n_l * impurity_l + n_r * impurity_r)

37                if ig > best_gain or best_threshold is None or best_feature is None:
38                    best_gain = ig
39                    best_feature = col
40                    best_threshold = threshold

41        self.feature = best_feature
42        self.gain = best_gain
43        self.threshold = best_threshold
44        if self.pre_pruning is False or self.depth < self.max_depth:
45            self._divide_tree(features, target, criterion)
46        else:
47            self.feature = None

48    def _divide_tree(self, features, target, criterion):
49        features_l = features[features[:, self.feature] <= self.threshold]
50        target_l = target[features[:, self.feature] <= self.threshold]
51        self.left = Tree(self.pre_pruning, self.max_depth)
52        self.left.depth = self.depth + 1
53        self.left.build(features_l, target_l, criterion)

54        features_r = features[features[:, self.feature] > self.threshold]
55        target_r = target[features[:, self.feature] > self.threshold]
56        self.right = Tree(self.pre_pruning, self.max_depth)
57        self.right.depth = self.depth + 1
58        self.right.build(features_r, target_r, criterion)

```

```

1
4 def _calc_impurity(self, criterion, target):
2     c = np.unique(target)
4     s = target.shape[0]
3
4     if criterion == 'gini':
4         return self._gini(target, c, s)
4     elif criterion == 'entropy':
5         return self._entropy(target, c, s)
4     elif criterion == 'error':
6         return self._classification_error(target, c, s)
4     elif criterion == 'mse':
7         return self._mse(target)
4     else:
8         return self._gini(target, c, s)
4
9 def _gini(self, target, classes, n_samples):
5     gini_index = 1.0
0     gini_index -= sum([(len(target[target==c]) / n_samples) ** 2 for c in classes])
5     return gini_index
1
5 def _entropy(self, target, classes, n_samples):
2     entropy = 0.0
5     for c in classes:
3         p = len(target[target==c]) / n_samples
5         if p > 0.0:
4             entropy -= p * np.log2(p)
5     return entropy
5
5 def _classification_error(self, target, classes, n_samples):
6     return 1.0 - max([len(target[target==c]) / n_samples for c in classes])
5
7 def _mse(self, target):
5     y_hat = np.mean(target)
8     return np.square(target - y_hat).mean()
5
9 # 決定木の事後剪定
6 def prune(self, method, max_depth, min_criterion, n_samples):
0     if self.feature is None:
6         return
1
6     self.left.prune(method, max_depth, min_criterion, n_samples)
2     self.right.prune(method, max_depth, min_criterion, n_samples)
6
3     pruning = False
6
4     if method == 'impurity' and self.left.feature is None and self.right.feature is None: # Leaf
6         if (self.gain * self.n_samples / n_samples) < min_criterion:
5             pruning = True
6     elif method == 'depth' and self.depth >= max_depth:
6         pruning = True
6
7     if pruning is True:
6         self.left = None
8         self.right = None
6         self.feature = None
9
7 def predict(self, d):
0     if self.feature is None: # Leaf
7         return self.label
1     else: # Node
7         if d[self.feature] <= self.threshold:
2             return self.left.predict(d)
7         else:
3             return self.right.predict(d)
7
4
7
5

```



1	
0	
7	
1	
0	
8	
1	
0	
9	
1	
1	
0	
1	
1	
1	
1	
1	
2	
1	
1	
3	
1	
1	
4	
1	
1	
5	
1	
1	
6	
1	
1	
7	
1	
1	
8	
1	
1	
9	
1	
2	
0	
1	
2	
1	
1	
2	
2	
1	
2	
3	
1	
2	
4	
1	
2	
5	
1	
2	
6	
1	
2	
7	
1	
2	
8	
1	
2	
9	

```

1
3
0
1
3
1
1
3
2
1
3
3
1
3
4
1
3
5
1
3
6
1
3
7

```

続いて、回帰木のための DecisionTreeRegressor クラスを書きます。

```

1  class DecisionTreeRegressor(object):
2      def __init__(self, criterion='mse', pre_pruning=False, pruning_method='depth', max_depth=3, min_
3  criterion=0.05):
4          self.root = None
5          self.criterion = criterion
6          self.pre_pruning = pre_pruning
7          self.pruning_method = pruning_method
8          self.max_depth = max_depth
9          self.min_criterion = min_criterion
10
11     def fit(self, features, target):
12         self.root = Tree(self.pre_pruning, self.max_depth)
13         self.root.build(features, target, self.criterion)
14         if self.pre_pruning is False: # post-pruning
15             self.root.prune(self.pruning_method, self.max_depth, self.min_criterion, self.root.n_sampl
16 es)
17
18     def predict(self, features):
19         return np.array([self.root.predict(f) for f in features])
20
21
22
23
24
25
26
27

```

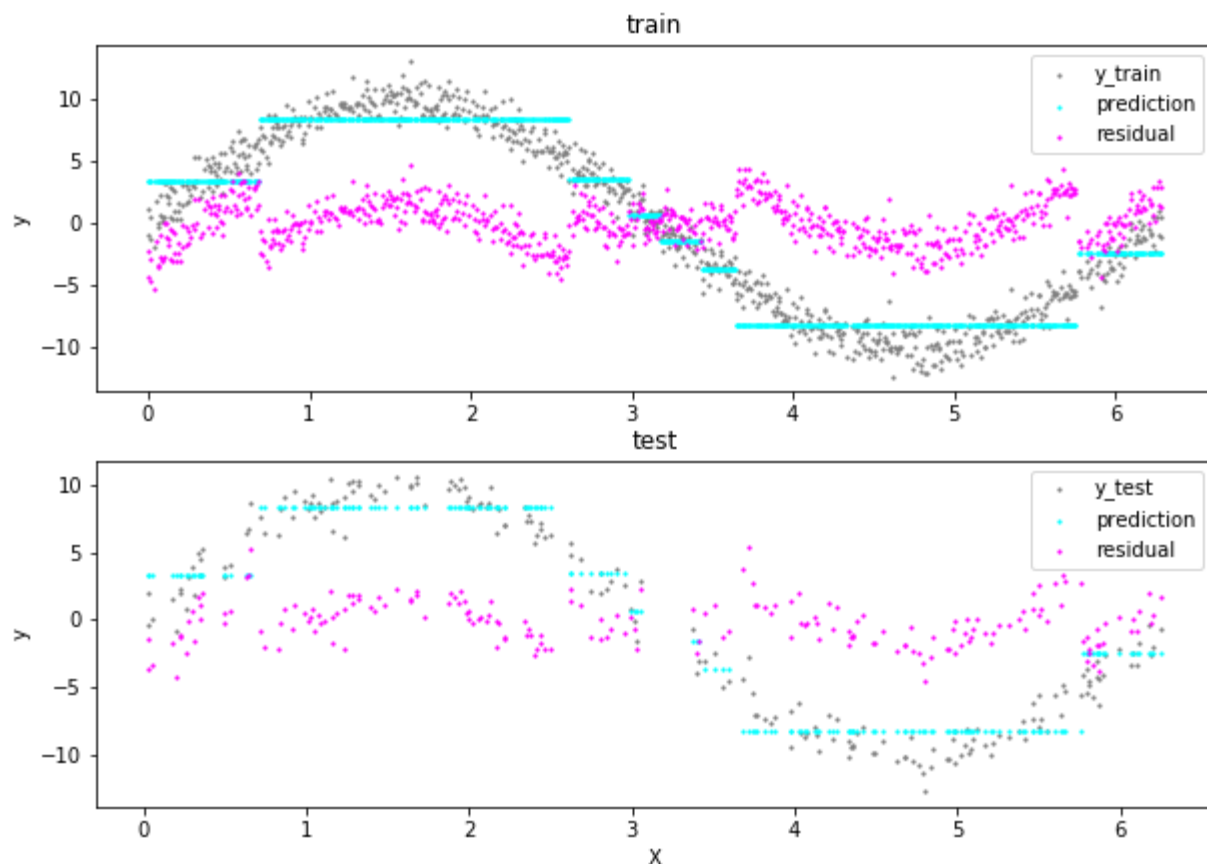
DecisionTreeRegressor クラスのインスタンスに対して fit() を呼び出し回帰木の分岐点を求めます。

```

1  X_train, X_test, y_train, y_test = generate_data()
2  regressor = DecisionTreeRegressor(criterion='mse', pre_pruning=True, pruning_method='depth', max_dept
3  h=3)
4  regressor.fit(X_train, y_train)

```

予測と残差を確認します。



```
1 def mse(y, pred):
2     return np.square(y - pred).mean()
3
4 print('MSE of the Train: %.2f, MSE of the Test: %.2f' % (mse(y_train, regressor.predict(X_train)), mse(y_test, regressor.predict(X_test))))
```

y には標準正規分布に従うノイズが付与されているため、当てはまりが十分である場合は MSE は概ね 1 程度となることが期待されます。しかし、訓練データの MSE は 2.92、テストデータの MSE は 2.86 となりました。従って、単一の最大深度3の回帰木では当てはまりが不十分であることがわかります。

GBDT

今回は Gradient Boosting の回帰アルゴリズムを提案した [Greedy Function Approximation: A Gradient Boosting Machine \[Jerome H. Friedman, 1999\]](#) 中の LS Boost を実装してみます。Gradient Boosting のアルゴリズムを上記の論文から引用します。

ALGORITHM 1 (Gradient_Boost).

```

1.  $F_0(\mathbf{x}) = \arg \min_{\rho} \sum_{i=1}^N L(y_i, \rho)$ 
2. For  $m = 1$  to  $M$  do:
3.  $\tilde{y}_i = -[\frac{\partial L(y_i, F(\mathbf{x}_i))}{\partial F(\mathbf{x}_i)}]_{F(\mathbf{x})=F_{m-1}(\mathbf{x})}, i = 1, N$ 
4.  $\mathbf{a}_m = \arg \min_{\mathbf{a}, \beta} \sum_{i=1}^N [\tilde{y}_i - \beta h(\mathbf{x}_i; \mathbf{a})]^2$ 
5.  $\rho_m = \arg \min_{\rho} \sum_{i=1}^N L(y_i, F_{m-1}(\mathbf{x}_i) + \rho h(\mathbf{x}_i; \mathbf{a}_m))$ 
6.  $F_m(\mathbf{x}) = F_{m-1}(\mathbf{x}) + \rho_m h(\mathbf{x}; \mathbf{a}_m)$ 
7. endFor
end Algorithm

```

数式の表記は以下となります。

- M : 弱学習器の数
- $L(y, F(x))$: 損失関数
- F_m : 加法モデル
- $h(x_i; a)$: 入力 x_i を引数にするパラメータ a を持つ弱学習器
- \tilde{y} : 損失関数の負の勾配
- ρ : 直線探索 (line search) によって与えられる係数 (勾配のステップサイズ)

GBDT では、最初のモデル F_0 を決めて、 $m=1$ から M まで損失関数が減少するように弱学習器を逐次的 (greedy stagewise) に最適化し、 F_m を更新します。

次に、GBDT の回帰アルゴリズムのひとつである LS_Boost (Algorithm 2) のアルゴリズムは以下です。

ALGORITHM 2 (LS_Boost).

```

 $F_0(\mathbf{x}) = \bar{y}$ 
For  $m = 1$  to  $M$  do:
   $\tilde{y}_i = y_i - F_{m-1}(\mathbf{x}_i), i = 1, N$ 
   $(\rho_m, \mathbf{a}_m) = \arg \min_{\mathbf{a}, \rho} \sum_{i=1}^N [\tilde{y}_i - \rho h(\mathbf{x}_i; \mathbf{a})]^2$ 
   $F_m(\mathbf{x}) = F_{m-1}(\mathbf{x}) + \rho_m h(\mathbf{x}; \mathbf{a}_m)$ 
endFor
end Algorithm

```

Gradient Boosting (Algorithm 1) における損失関数を二乗誤差, $L(y, F(x)) = (y - F(x))^2 / 2$ とすると Algorithm 1 の3行目の \tilde{y} はシンプルに残差 $\tilde{y} = y_i - F_{m-1}(x_i)$ となります。また、 $\rho_m = \beta_m$ となるため ρ は α と共に最適化できます。ちなみに、GBDT の回帰アルゴリズムを GBRT (Gradient Boosted Regression Trees) と表している文献も多くあります。

早速、LS_Boost を実装し実験してみます。弱学習器である回帰木の最大深度は先程と同様に 3 としています。

```

1 X_train, X_test, y_train, y_test = generate_data()
2
3 M = 10
4 predictions_history = [np.repeat(y_train.mean(), len(y_train))]
5 test_predictions_history = [np.repeat(y_test.mean(), len(y_test))]
6

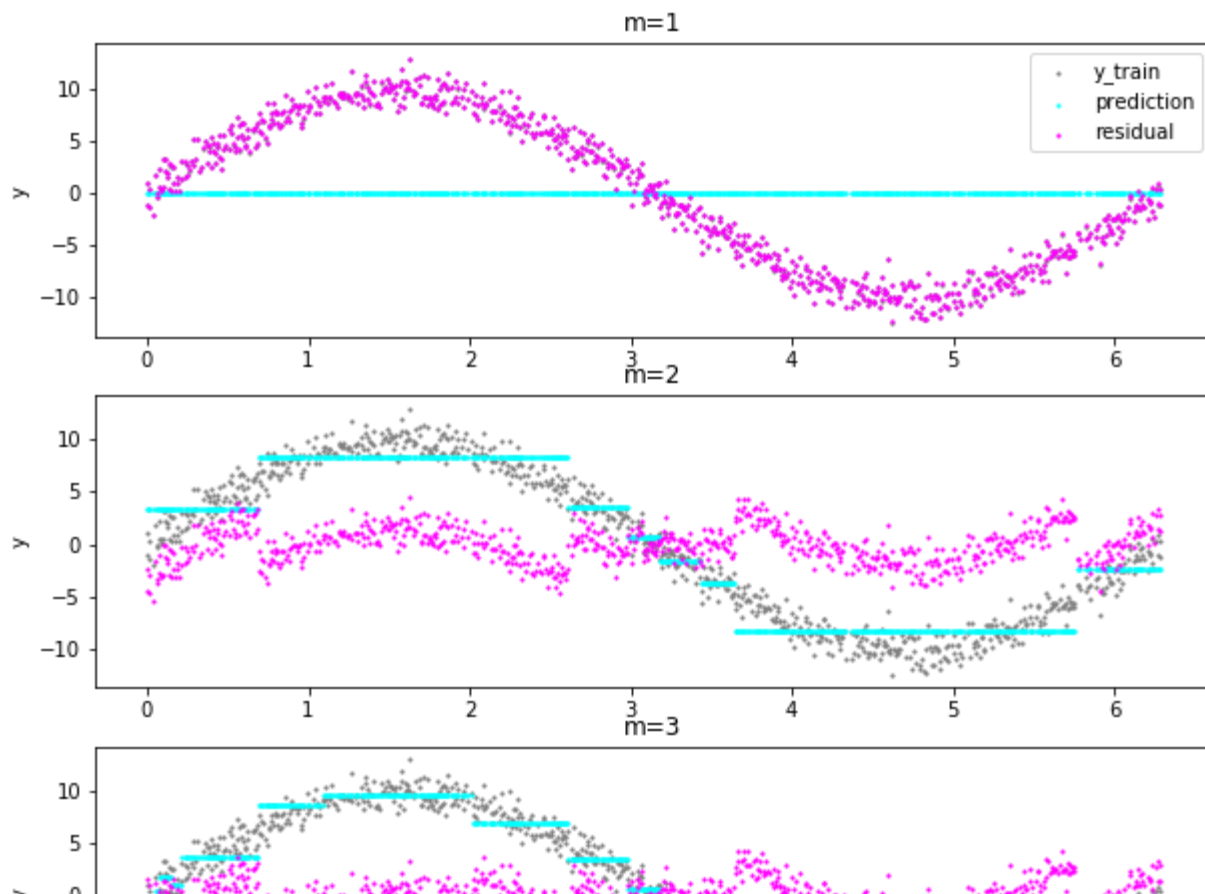
```

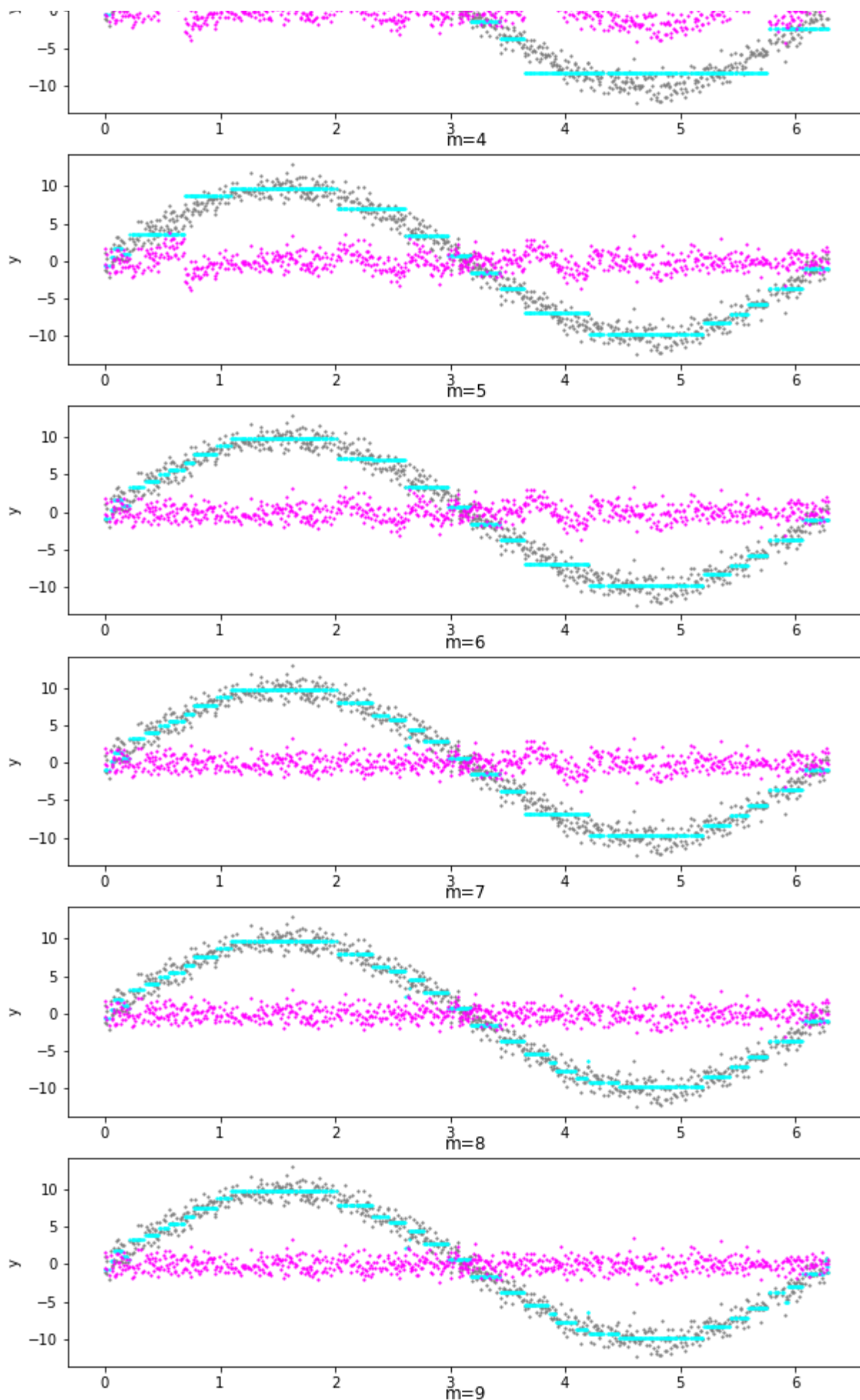
```

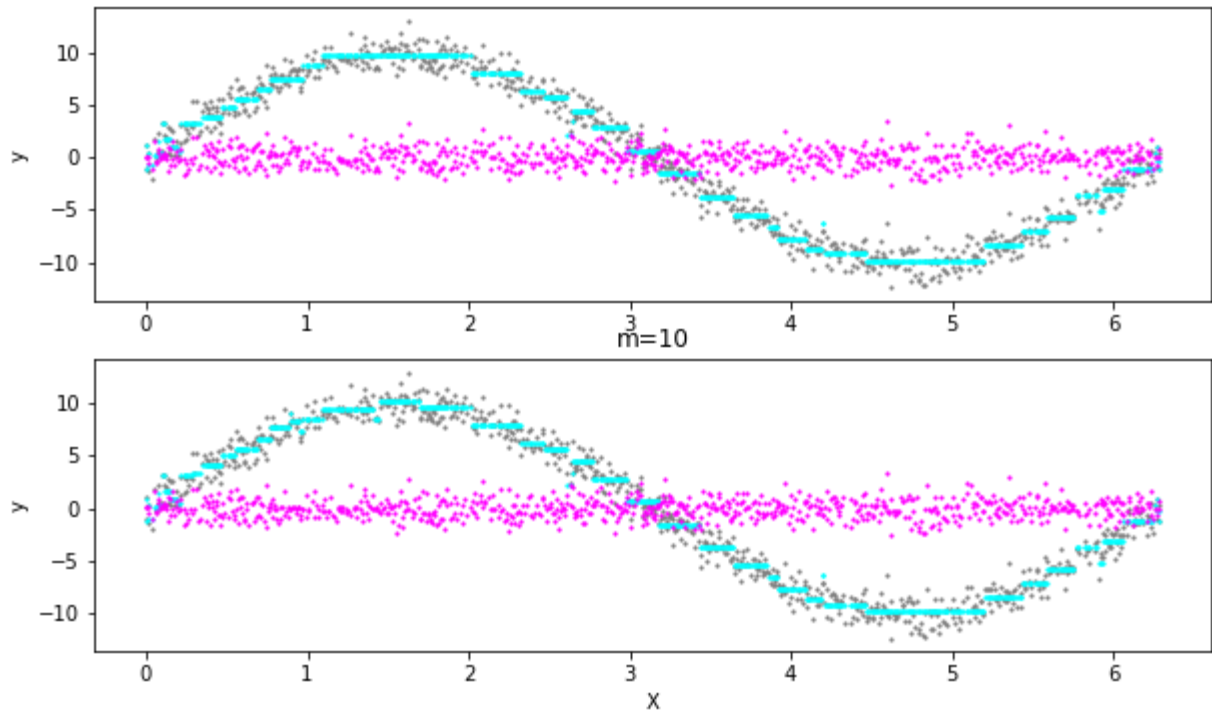
7 # LS_TreeBoost (Algorithm-2)
8 for m in range(M):
9     y_tilde = y_train - predictions_history[-1]
10     base_learner = DecisionTreeRegressor(criterion='mse', pre_pruning=True, pruning_method='depth', m
ax_depth=3)
11     base_learner.fit(X_train, y_tilde)
12
13     prediction = predictions_history[-1] + base_learner.predict(X_train)
14     test_prediction = test_predictions_history[-1] + base_learner.predict(X_test)
15
16     train_mse = mse(y_train, prediction)
17     test_mse = mse(y_test, test_prediction)
18
19     predictions_history.append(prediction)
20     test_predictions_history.append(test_prediction)
21
22     print("[%d] column: %d, threshold: %f, mse-of-train: %.2f, mse-of-test: %.2f\n" %
23           (m+1, base_learner.root.feature, base_learner.root.threshold,
24            train_mse, test_mse) + "-" * 50)
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99

```

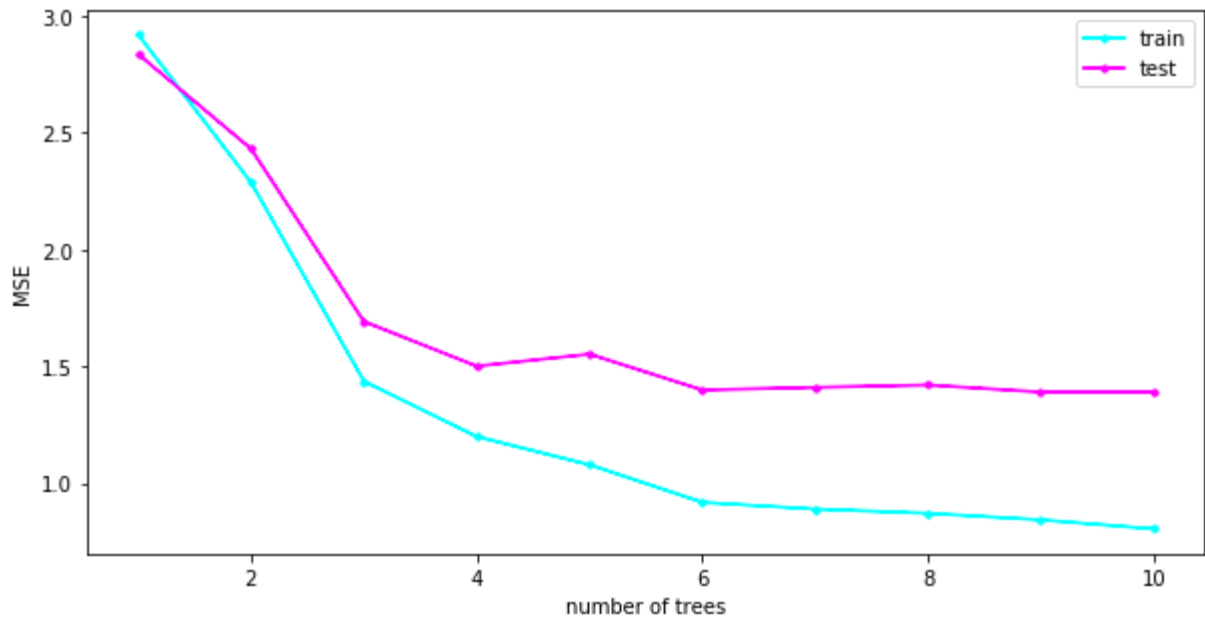
残差に対して fitting を行う回帰木を逐次的に追加しモデル F を更新していく様子を確認します。







m=1 から 10 までの MSE の推移は以下です。



m=10 の段階で、訓練データの MSE は 0.81、テストデータの MSE は 1.39 となりました。一方で、m=8 前後からノイズに fitting している箇所が出てきています。論文中では、GBDT の正則化パラメータである弱学習器の数 M や学習率 η により過学習を制御する方法が紹介されています。

おわりに

この記事では、GBDT の回帰アルゴリズムのひとつである LS Boost を実装し、LS Boost が残差に対して fitting を行う回帰木を逐次的に追加しモデルを更新していく様子を確認しました。今回の実装は GBDT のアルゴリズムを理解するためのものでしたが、Kaggle に代表されるデータサイエンスコンペティションで人気を集めている XGBoost や LightGBM では GBDT を大規模データに適用するための様々な高速化・効率化の手法が実装されています。[1,2]

参考文献

- [1] XGBoost: A Scalable Tree Boosting System [Chen, 2016]
- [2] LightGBM: A Highly Efficient Gradient Boosting Decision Tree [Ke et al., 2017]
- [3] Python機械学習プログラミング 達人データサイエンティストによる理論と実践

記事の感想やリクエストはありませんか？

送信

ツイート

O

#AI#Python#機械学習

2016年11月、データサイエンティストとして中途入社。時系列分析や異常検知、情報推薦に特に興味があります。クロスバイク、映画鑑賞、猫が好き。

Recommends

こちらもおすすめ

2018.11.9

Data Science

S

Tech

2020.5.18 GCP2020.5

GCP2020.3

ISMS認証