

Boosting Methods

MODULES

DATE: _____

Section Overview:

- Boosting and Meta-learning
- AdaBoost (Adaptive Boosting) Theory
- Example of AdaBoost
- Gradient Boosting Theory
- Example of Gradient Boosting

2 - Boosting Motivation and History

The concept of Boosting is not actually a machine learning algorithm, it is methodology applied to an existing machine learning algorithm, most commonly applied to decision tree.

- main formula for boosting

$$F_T(x) = \sum_{t=1}^T f_t(x)$$

$$f_t(x) = \alpha_t h_t(x) \quad \downarrow$$

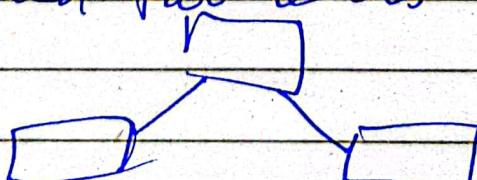
- Implies that a combination of estimators with an applied coefficient could act as an effective ensemble estimator.

- Note $h(x)$ can in theory be any machine learning algorithm (estimator/learner)
- Can an ensemble of weaker learners (very simple models) be a strong learner when combined?
- For decision tree models, we can use simple tree in place of $h(x)$ and combine them with the coefficients on each model.
- The idea of gradient boosting originated from Leo Breiman when he observed that boosting can be interpreted as an optimized ~~function~~ algorithm on a cost function in publications in the late 1990s.
 - later on Jerome H Friedman and many others developed more explicit formulations of gradient boosting.
- Also in the late 1990s Yoav Freund and Robert Schapire developed the AdaBoost (Adaptive Boosting) algorithm, which also combines weak learner in an ensemble to create a ~~sim~~ stronger model.

- Let's continue by focusing first on AdaBoost and building an understanding of how to combine weak learners to create a strong estimator.
- We will also explore why Decision Trees are so suited for boosting

3. AdaBoost Theory and Intuition

- AdaBoost (Adaptive Boosting) work by using an ensemble of weak learners and then combining them through the use of weighted sum
- AdaBoost adapts by using previously created weak learners in order to adjust misclassified instances for the next created weak learner.
- What is a weak learner?
 - A weak learner is a model that is too simple to perform well on its own -
- The weakest decision tree possible would be a stump | one node and two leaves -



M T W T F S

DATE:

Slides have the necessary data
Check them

004 : AdaBoost Coding part 1

Coding

```
. df = pd.read_csv('.../Data/mushrooms.csv')  
df.head()
```

```
sns.countplot(data=df, x='class')  
df.describe()
```

```
df.describe().transpose().reset_index()
```

```
feat = df.sort_values('unique')
```

```
sns.barplot(data=feat, x='index', y='unique')  
plt.xticks(rotation=90);
```

```
x = df.drop(['class'], axis=1)  
x.isnull().sum()
```

```
x = pd.getdummies(x, drop_first=True)
```

```
y = df['class']
```

`train-test-split(X, Y)`

Part 2 Adaboost Classifier

```
from sklearn.ensemble import AdaBoostClassifier
model = AdaBoostClassifier(n_estimators=1)
model.fit(X-train, y-train)
```

metrics

```
predictions = model.predict(X-test)
```

```
predictions
print(classification_report)
```

```
model.feature_importance_.argmax()
```

```
X.columns[22]
```

$\approx 10 \text{ or } 11$

```
sns.countplot(data=df, hue="class")
```

```
len(X.columns)
```

for loop n in range(1, 96):

```
model = AdaBoostClassifier(n_estimators=n)
model.fit(X-train, y-train)
```

```
preds = model.predict(X-test)
```

err = 1 - Accuracy(y-test, y-pred)

error_rates.append(err)

```
plt.plot(range(1, 96), error_rates)
```

$\text{fearts} = \text{prod}(\text{Data}[\text{pm} \geq \text{val}])$

$\text{model} =$
 model_
 Input

Gradient Boosting

Theory

Gradient boosting is very similar idea to AdaBoost where weak learners are created in series in order to produce a strong ensemble method-model.

- Gradient Boosting make use of residual error for learning.
- Gradient Boosting vs AdaBoost:
 - Larger tree allowed in Gradient Boosting
 - Learning rate coefficient same for all weak learners.
- Gradient series learning is based on training on the residual of previous model.

y	\hat{y}	e	f_1	$F_1 = \hat{y} + f_1$
5000000	504,000	9000	8000	501,000
462000	509,000	47000	50000	459000
565000	509,000	56000	50000	559,000

- Create new model to predict the error
- update prediction using error prediction
- Gradient Boosting Process:-

$$F_m = F_{m-1} + f_m$$

$$F_m = F_{m-1} + (\text{learning rate} \cdot f_m)$$

- create initial model
- Train another model on error

$$\Rightarrow e = y - f_0$$

- create new prediction

$$F_1 = f_0 + \eta f_1$$

- Repeat as needed

$$F_m = f_{m-1} + \eta f_m$$

- Note, for classification we can use the logit as an error metric:

$$\hat{y} = \log\left(\frac{\hat{p}}{1-\hat{p}}\right) \quad \hat{p} = \frac{1}{1+e^{-\theta}}$$

M T W T F S

DATE: _____

- Note: the learning rate is not same for each new model in the series, it is not unique to each subsequent model (unlike AdaBoost's alpha coefficient).
- Gradient Boosting is fairly robust to overfitting, allowing for the number of estimators to be set high default ($t=100$)
- Gradient Boosting Intuition:
 - we optimize the series of trees by learning on the residuals, forcing subsequent trees to attempt to correct for the error in the previous trees.
 - The trade off is training time,
 - A learning rate is between 0-1, which means a very low value would mean each subsequent tree has little "say", meaning more trees need to be created, causing a larger computation training time.

Gradient Boosting Coding

NOTES

DATE: _____

Same as beta before for Adaboost
Input

from sklearn.ensemble.GradientBoostingClassifier

param_grid = {
 'n_estimators': [50, 100],
 'learning_rate': [0.1, 0.05, 0.2],
 'max_depth': [3, 4, 5]}

gb_model = GradientBoostingClassifier()

grid = GridSearchCV(gb_model,
 param_grid).

grid.fit(X_train, y_train)

metrics