



مبانی هوش محاسباتی (پروژه سوم)

کیارش آستان بوس ۴۰۰۱۲۶۲۵۷۰

محمد رضا نادری ۴۰۰۱۲۶۲۳۸۶

مهر ۱۴۰۳

فهرست مطالب

۲	Visualize (۱)
۳	Bagging (۲)
۳ Training ۱.(۲)
۳ Accuracy by n_estimators ۲.(۲)
۴ Decision Boundaries ۳.(۲)
۴ Model Evaluate ۴.(۲)
۵ Accuracy by n_estimators (Library) ۵.(۲)
۶ Decision Boundaries (Library) ۶.(۲)
۶ Model Evaluate (Library) ۷.(۲)
۷	RandomForest (۳)
۷ Training ۱.(۳)
۷ Accuracy by n_estimators ۲.(۳)
۸ Decision Boundaries ۳.(۳)
۸ Model Evaluate ۴.(۳)
۹ Accuracy by n_estimators (Library) ۵.(۳)
۱۰ Decision Boundaries (Library) ۶.(۳)
۱۰ Model Evaluate (Library) ۷.(۳)
۱۲	AdaBoost (۴)
۱۲ Training ۱.(۴)
۱۲ Accuracy by n_estimators ۲.(۴)
۱۳ Decision Boundaries ۳.(۴)
۱۴ Model Evaluate ۴.(۴)
۱۵	Stacked Learners (۵)
۱۵ Training ۱.(۵)
۱۵ Model Evaluate ۲.(۵)
۱۶ (Extra) overfitting ۳.(۵)

Visualize (۱)



شکل ۱:

با نمایش داده ها میتوان دید دو کلاس داریم که هر کدام به دو بخش تقسیم شده اند. همچنین تراکم نیز در دو دسته سمت چپ داده ها بیشتر است. در هر دسته نیز تعداد اندکی از کلاس مقابل وجود دارد که نویز محسوب می شود. و همچنین چندین داده پرت از هر دو کلاس داریم که به تراکمی نزدیک نیستند.

Bagging (۲)

Training ۱.(۲)

: GridSearch Parameters

```
model = BaggingClassifierManual()
param_grid = {
    "n_estimators": [10, 15, 25, 50, 100],
    "max_depth": [5, 7, 10]
}
```

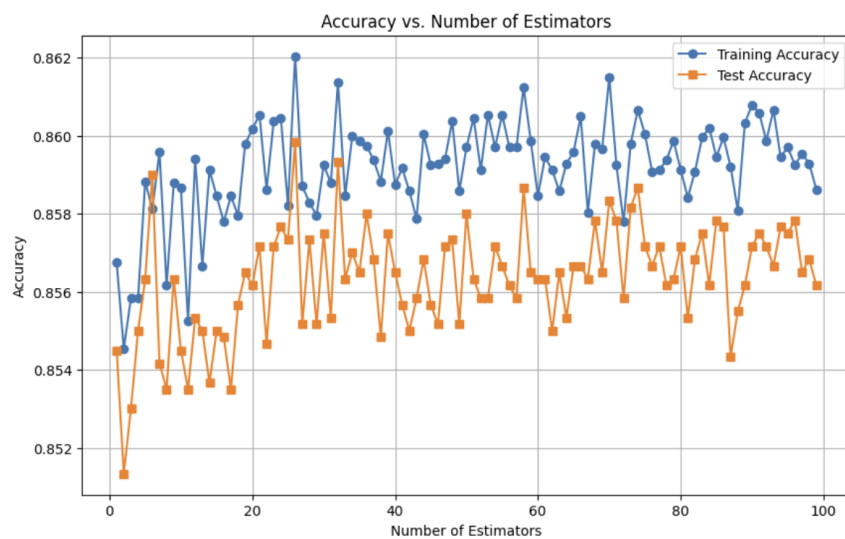
شکل ۲:

: Best Parameters

• $n_{estimators}$: ۱۰

• max_depth : ۵

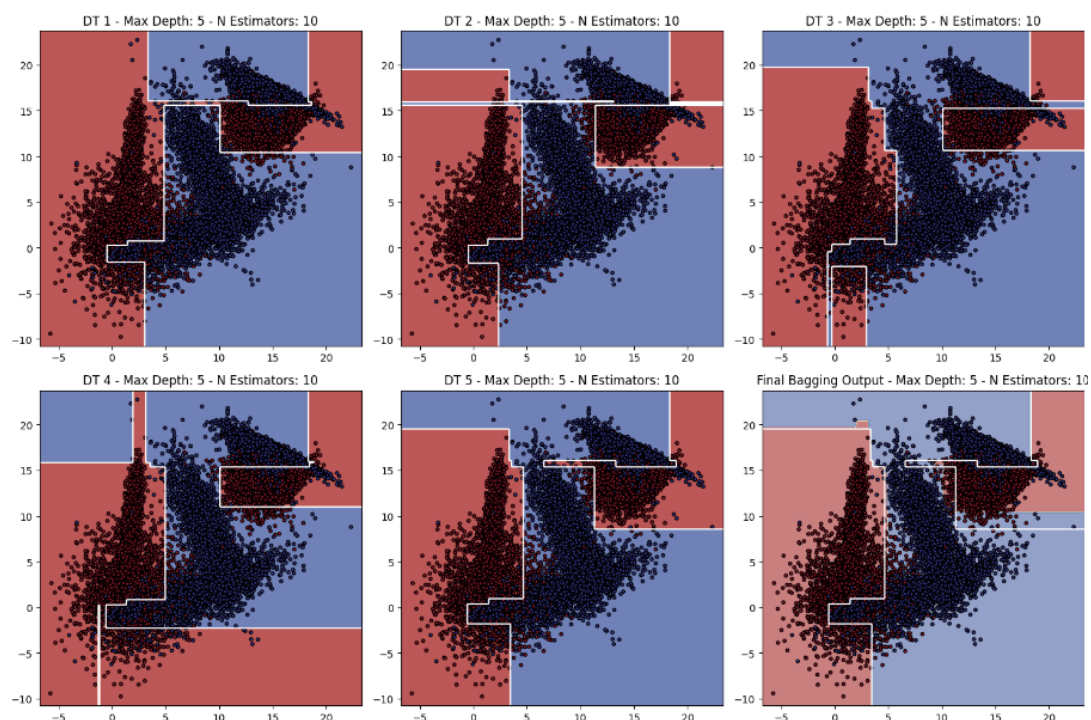
۲.(۲) Accuracy by $n_{estimators}$



شکل ۳:

با توجه به نمودار بالا که accuracy را فقط با توجه به پارامتر Estimators سنجیده است میتوان به نتیجه رسید که افزایش این پارامتر میتواند تا حدودی مدل را بهتر کند اما صرفاً مدل به این پارامتر بستگی ندارد و باید به پارامتر دیگری مثل عمق درخت هم توجه کرد

۳.(۲) Decision Boundaries



شکل ۴:

هر زیرنمودار نشان‌دهنده یکی از مدل‌های پایه (درخت تصمیم) است که با استفاده از داده‌های نمونه‌برداری‌شده‌ی متفاوت (Bootstrap) و با حداکثر عمق (Depth Max) مشخص آموزش دیده است. این درخت‌ها به صورت مستقل از هم آموزش داده می‌شوند.

هر مدل یک مرز تصمیم متفاوت ایجاد کرده است که به طور کامل از دیگر مدل‌ها متمایز است. این تفاوت به دلیل استفاده از زیرمجموعه‌های تصادفی از داده‌ها و ویژگی‌ها است.

آخرین زیرنمودار، نتیجه‌ی نهایی الگوریتم Bagging است. در این نمودار می‌توان مشاهده کرد که مرزهای تصمیم نهایی بسیار هموارتر و دقیق‌تر هستند. ترکیب مدل‌ها باعث شده که خطاهای هر مدل به طور موثری کاهش یافته و عملکرد کلی بهبود پیدا کند.

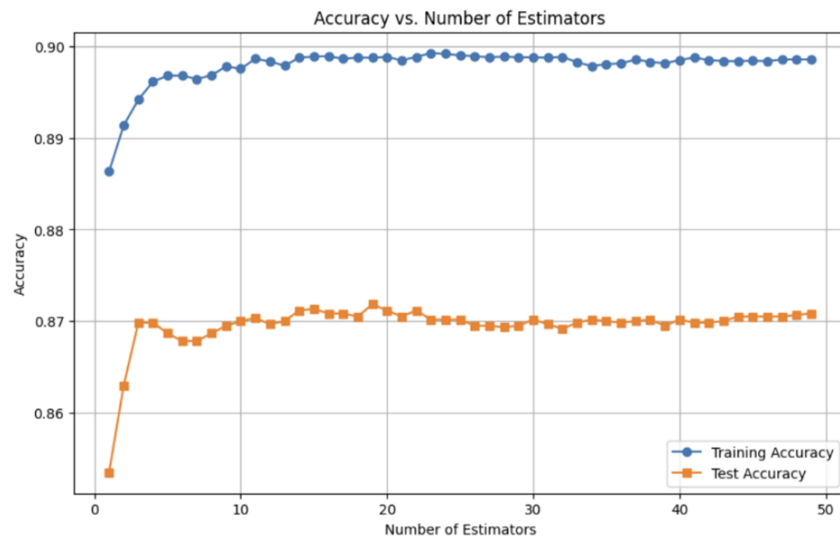
۴.(۲) Model Evaluate

	Training	Test
Accuracy	0.859167	0.858167
Precision	0.859447	0.858291
Recall	0.859167	0.858167
F1 Score	0.859142	0.858146

شکل ۵:

نزدیک بود accuracy در داده های train و test نشان میدهد که مدل خوبی را در خروجی داریم و overfit هم رخ نداده است

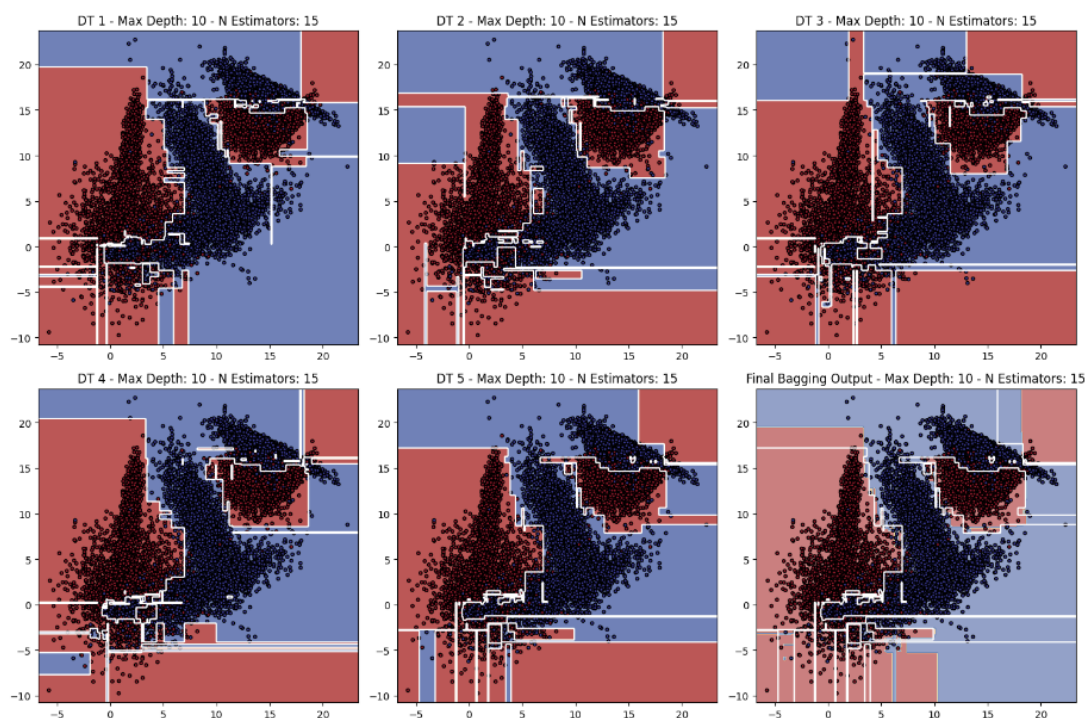
۵.۲) Accuracy by n_estimators (Library)



شکل ۶:

در این نمودار میبینیم که رشد از تعداد حدودا ۱۵ متوقف شده است و رشد چندان زیادی دیده نمیشود. یکی از دلایل تفاوت در این نمودار و نمودار قبلی میتواند در این باشد که در اینجا ما از تابع آماده استفاده کرده ایم که میتواند بسیار بهینه تر و بهتر پیاده سازی شده باشد.

Decision Boundaries (Library) ۶.۲



شکل ۷:

در اینجا هم هر زیرنمودار نشان‌دهنده یکی از مدل‌های پایه (درخت تصمیم) است که با استفاده از داده‌های نمونه‌برداری‌شده‌ی متفاوت (Bootstrap) و با حداکثر عمق (Depth Max) مشخص آموزش دیده است. این درخت‌ها به صورت مستقل از هم آموزش داده می‌شوند. آخرین زیرنمودار، نتیجه‌ی نهایی الگوریتم Bagging است. چندان تفاوتی با کلاس بندی تابع دستی ندارد

Model Evaluate (Library) ۷.۲

	Training	Test
Accuracy	0.898875	0.871333
Precision	0.899269	0.871802
Recall	0.898875	0.871333
F1 Score	0.898848	0.871306

شکل ۸:

این ارزیابی نشان می‌دهد که تفاوت چندان زیادی بین تابعی که به صورت دستی پیاده سازی شده و تابع آماده در کتابخانه sklearn نیست و نتیجه خوبی به ما داده است

RandomForest (۳)

Training ۱.(۳)

: GridSearch Parameters

```
model = RandomForestManual()
param_grid = {
    "n_estimators": [10, 15],
    "max_depth": [5, 7, None],
    "max_features": [None, 'sqrt', 'log2']
}
```

شکل ۹:

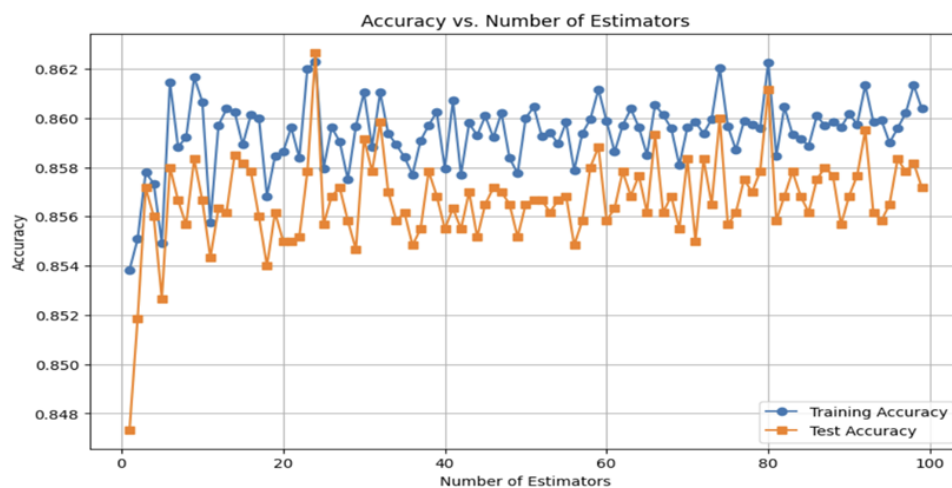
: Best Parameters

• n_estimators : ۱۰

• max_depth : ۱۵

• max_features : None

Accuracy by n_estimators ۲.(۳)

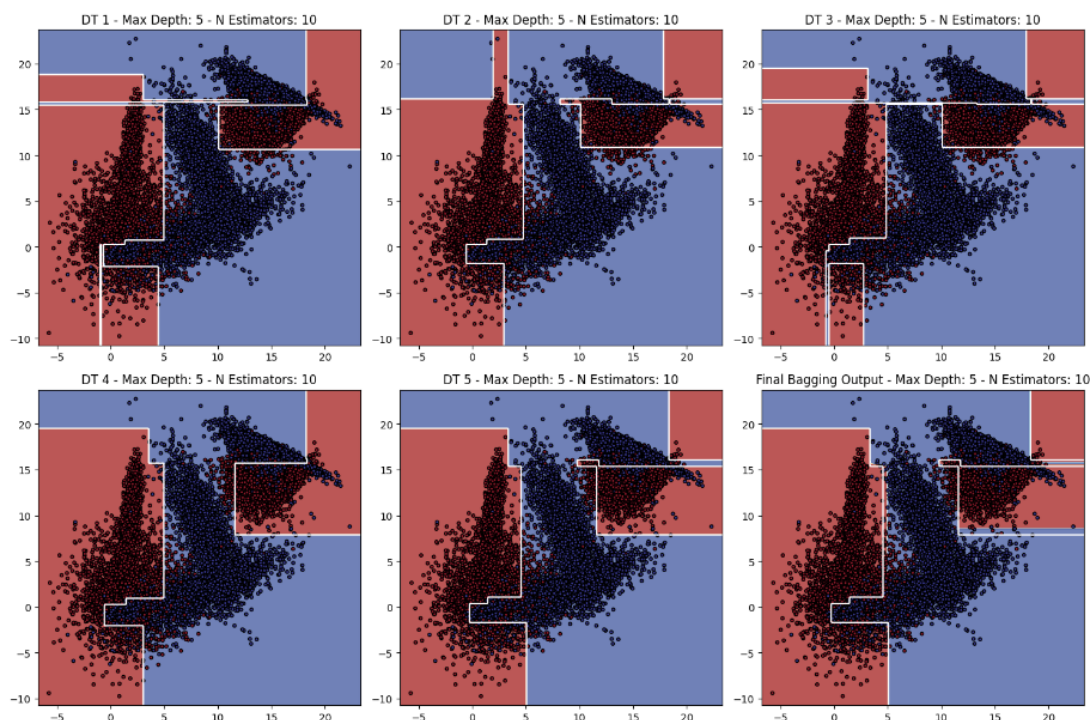


شکل ۱۰:

با توجه به نمودار بالا که accuracy را فقط با توجه به پارامتر Estimators سنجیده است میتوان به نتیجه رسید که افزایش این پارامتر میتواند تا حدودی مدل را بهتر کند اما صرفاً مدل به این پارامتر بستگی ندارد و باید به پارامتر دیگری مثل عمق درخت هم توجه کرد

در این نمودار دیده میشود که در یکی از مقادیر accuracy داده تست بیشتر شده است این به دلیل بهینه نبودن تابع میتواند باشد که منجر به این اتفاق شده است

۳.(۳) Decision Boundaries



شکل ۱۱:

در ۵ نمودار اول هر درخت تصمیم‌گیری (DT) تنها بخشی از داده‌ها را از طریق Sampling Bootstrap دریافت می‌کند و بر اساس زیرمجموعه‌ای از ویژگی‌ها آموزش می‌بیند. مرزهای تصمیم درخت‌ها به صورت تکه تکه هستند، زیرا:

- هر درخت فقط بر اساس داده‌های محدود آموزش داده شده است.
 - محدودیت در max_depth باعث می‌شود مدل بیش از حد پیچیده نشود و مرزهای ساده‌تری ارائه کند.
 - این مرزها ممکن است به تنهایی دقیق نباشند و منجر به پیش‌بینی‌های نادرست شوند.
- نمودار آخر (پایین سمت راست) نشان‌دهنده خروجی نهایی الگوریتم Random Forest است. این خروجی با استفاده از رأی‌گیری اکثریت (majority voting) بین پیش‌بینی‌های تمام درخت‌ها به دست آمده است.

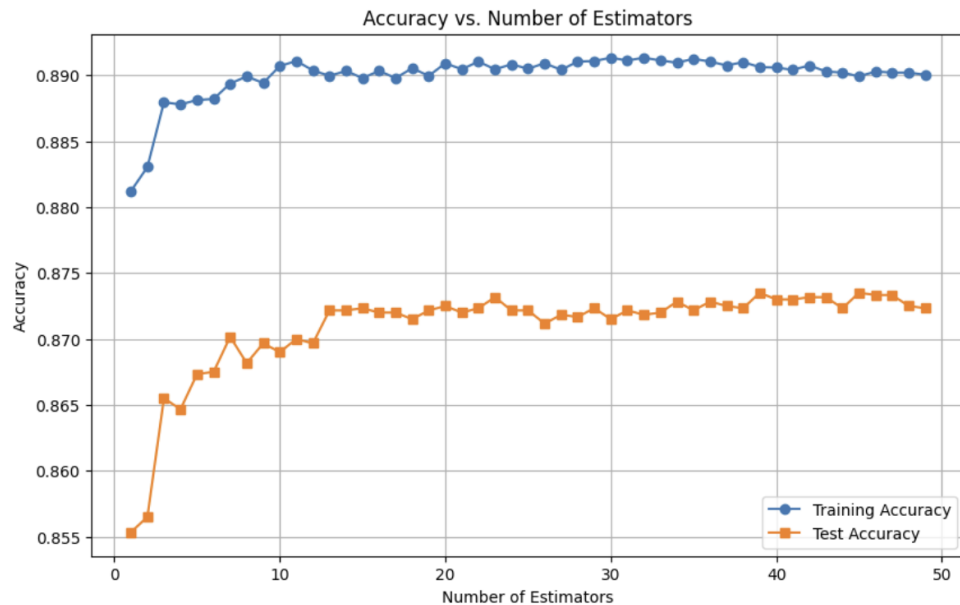
۴.(۳) Model Evaluate

	Training	Test
Accuracy	0.857000	0.854000
Precision	0.857324	0.854176
Recall	0.857000	0.854000
F1 Score	0.856971	0.853972

شکل ۱۲:

ارزیابی این مدل با پارامترهای انتخاب شده به شرح زیر است که برای دیتا درصد نسبتاً خوبی را داده است و overfit رخ نداده است.

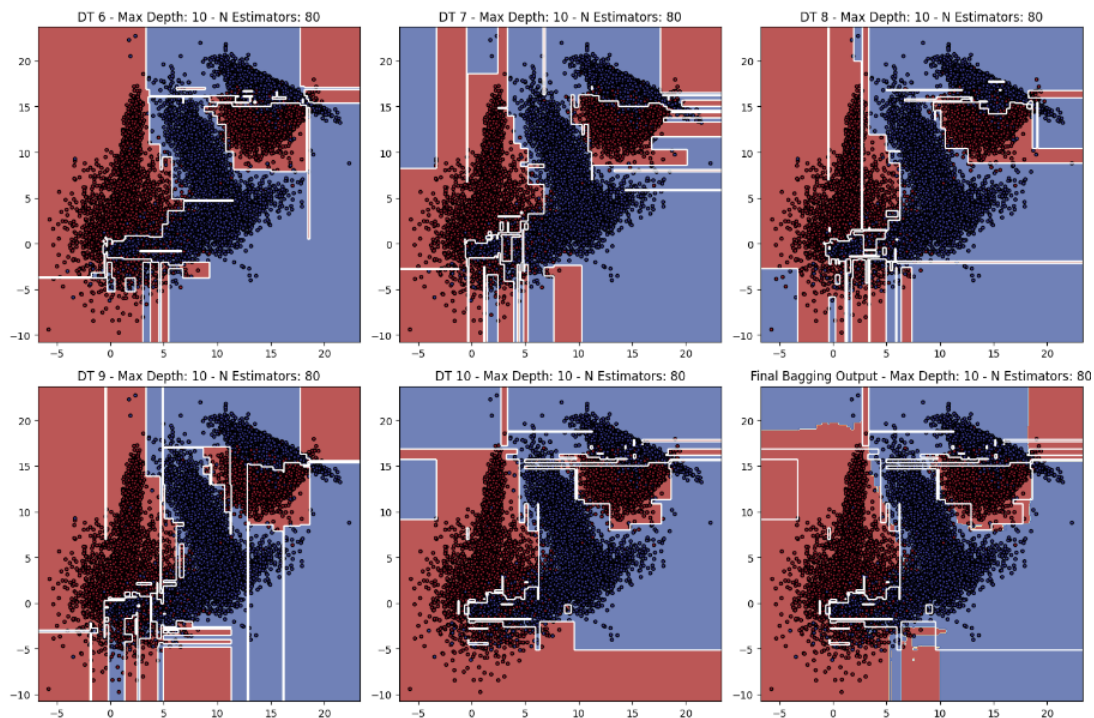
۵.۳ Accuracy by n_estimators (Library)



شکل ۱۳:

در این نمودار میبینیم که رشد از تعداد حدوداً ۱۳ متوقف شده است و رشد چندان زیادی دیده نمیشود و حتی در مقدار های ۴۰ به بعد اندکی کاهش پیدا کرده است. یکی از دلایل تفاوت در این نمودار و نمودار قبلی میتواند در این باشد که در اینجا ما از تابع آماده استفاده کرده ایم که میتواند بسیار بهینه تر و بهتر پیاده سازی شده باشد.

Decision Boundaries (Library) ۶.(۳)



شکل ۱۴:

این درخت‌ها به دلیل تکنیک‌های زیرمجموعه سازی ویژگی‌ها و داده‌ها که در جنگل تصادفی استفاده می‌شوند، پیش‌بینی‌های متفاوتی ارائه می‌دهند.

اختلاف در این گراف‌ها نشان‌دهنده این است که هر درخت به صورت جزئی داده‌ها را مدل می‌کند. در گراف‌های DT1 تا DT5، مرزهای تصمیم اغلب ناپایدار و غیر دقیق هستند. این رفتار به دلیل پیچیدگی محدود هر درخت تصمیم و عدم تعمیم‌پذیری مناسب یک درخت است. در گراف نهایی، مرزهای تصمیم بهبود یافته و صاف‌تر شده‌اند، زیرا ترکیب درخت‌ها نویز را کاهش داده و منجر به مدل قوی‌تری شده است.

Model Evaluate (Library) ۷.(۳)

	Metric	Train	Test
0	Accuracy	0.894542	0.871500
1	Precision	0.895000	0.872111
2	Recall	0.894542	0.871500
3	F1 Score	0.894508	0.871462

شکل ۱۵:

حدود دو درصد در تست و ۴ درصد در تمرین بهبود داشته است و در کل دقت قابل قبول و تعمیم پذیری بالایی را نتیجه می‌دهد.

AdaBoost (۴)

Training ۱.(۴)

: GridSearch Parameters

```
param_grid = {
    "n_estimators": [50, 100, 200],
    "learning_rate": [0.01, 0.1, 1],
    'estimator': [DecisionTreeClassifier(max_depth=1),
                  DecisionTreeClassifier(max_depth=2),
                  DecisionTreeClassifier(max_depth=3),
                  DecisionTreeClassifier(max_depth=4),
                  DecisionTreeClassifier(max_depth=5)]
}
```

شکل ۱۶:

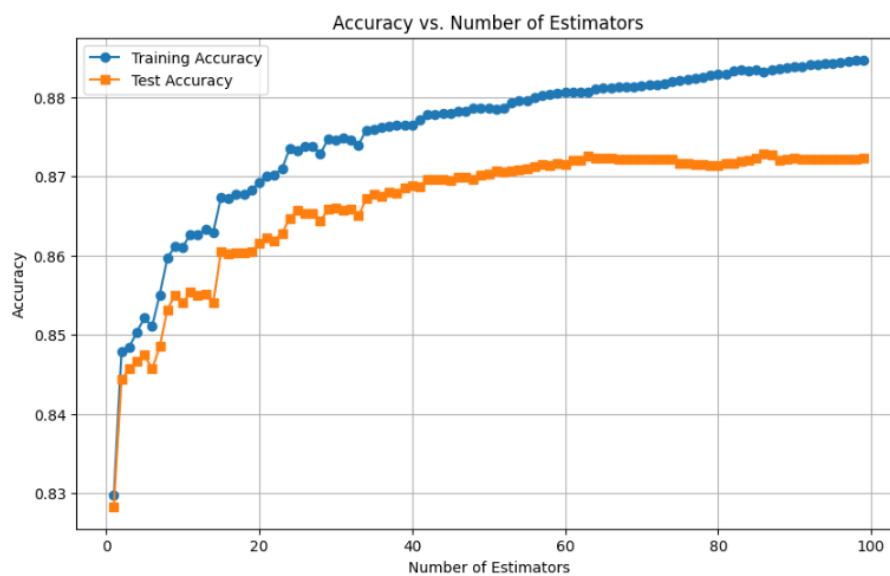
: Best Parameters

• n_estimators : ۱۰۰

• learning_rate : ۰.۱

• estimator: DecisionTreeClassifier(max_depth=۴)

Accuracy by n_estimators ۲.(۴)



شکل ۱۷:

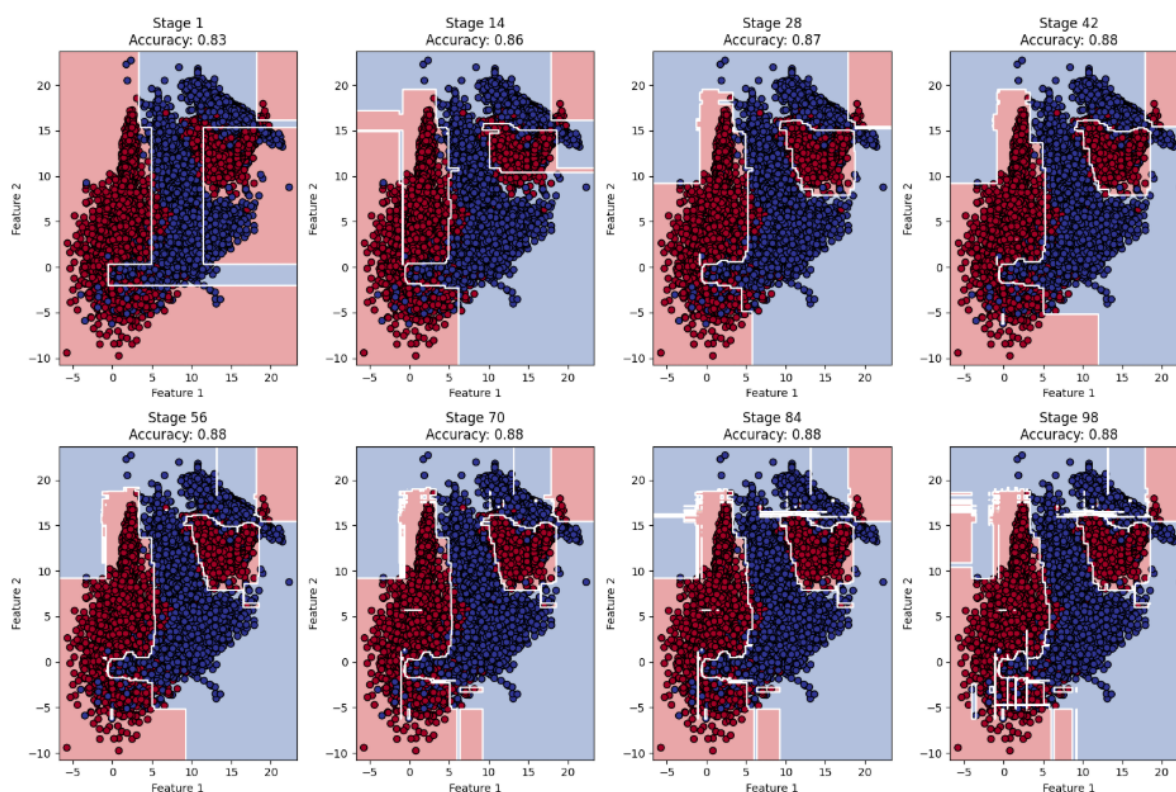
: Training Data

با افزایش تعداد تخمین زن ها، خطای داده‌های آموزشی را به صورت تکراری کاهش دهد. الگوریتم به تطبیق با داده‌های آموزشی ادامه می‌دهد و باعث افزایش دقت روی مجموعه آموزشی می‌شود. اما شیب افزایش دقت به مرور کاهش می‌یابد چون الگوریتم با داده‌های آموزشی تطبیق پیدا می‌کند.

: Test Data

دقت داده‌های تست در ابتدا افزایش می‌یابد، زیرا مدل با تعداد تخمین زن‌های بیشتر بهتر تعمیم می‌دهد. اما پس از یک نقطه مشخص (حدود ۶۰ تخمین زن)، شیب تقریباً صفر می‌شود، که نشان می‌دهد اضافه کردن تخمین زن‌های بیشتر تأثیر قابل توجهی در بهبود تعمیم ندارد. این مسطح شدن به این دلیل اتفاق می‌افتد که AdaBoost ممکن است با افزایش بیش از حد تعداد تخمین زن‌ها، بیش از حد به داده‌های آموزشی تطبیق پیدا کند و پیچیدگی اضافی دیگر به پیش‌بینی بهتر داده‌های جدید کمکی نمی‌کند.

۳.۴ Decision Boundaries



شکل ۱۸:

در این قسمت می‌توانیم مرز بین طبقه‌بندی هارا مشاهده کنیم. همانطور که در تصاویر مشخص است، هرچقدر به پایان الگوریتم نزدیک تر می‌شویم دقت افزایش می‌یابد و از دقت ۸۳ درصد به ۸۷ درصد رسیده‌ایم.

همچنین در بهبود مرزهای تصمیم‌گیری به عنوان مثال بین مرحله ۱ و ۱۴، داده‌های سمت راست کلاس آبی، بخشی زیادی از آنها در ابتدا به کلاس قرمز پیش‌بینی می‌شده‌اند که در مرحله ۱۴ آبی پیش‌بینی شده‌اند. اما با رفع این مشکل بخشی از داده‌های قرمز (بین ۱۰ و ۲۰ افقی و ۱۰ و ۸ عمودی) به کلاس آبی پیش‌بینی شده‌اند که در مرحله ۲۸ این مشکل رفع شده.

هرچقدر مرحله جلوتر می‌رود مدل داده‌های جزئی که ممکن است از تراکم کلاس خودشان دور باشد هم درست پیش‌بینی میکند مثلاً دو داده کلاس آبی در مست چپ (بین ۵- و ۵+ افقی و ۸- و ۲- عمودی) که تا قبل مرحله ۹۸ برای کلاس قرمز پیش‌بینی می‌شده‌اند نیز به درستی پیش‌بینی خواهند شد. حتی نواحی که در آنها داده‌ای نیست نیز تقسیم بندی شده‌اند (بین ۵ و ۱۰ افقی و ۱۰- و ۵- عمودی) یا بخش بزرگی از دیتا (بین ۱۰- و ۵- عمودی) در ابتدا کاملاً کلاس قرمز پیش‌بینی می‌شد و به مرور حدوداً این ناحیه نصف برای کلاس آبی شد و نصف برای کلاس قرمز.

در کل میتواند دید هرچه جلوتر می‌رویم نواحی کلاس‌ها از مربعی بودن تغییر میکند و شکل دیتا را به خود می‌گیرد مخصوصاً برای کلاس قرمز (بین ۱۰ و ۲۰ افقی و ۱۰ و ۱۵ عمودی) که کاملاً شکل کلاس قرمز را به خود گرفته.

Model Evaluate ۴.(۴)

	Training	Test
Accuracy	0.884708	0.872333
Precision	0.884827	0.872589
Recall	0.884708	0.872333
F1 Score	0.884698	0.872321

شکل ۱۹:

تفاوت بین دقت داده‌های تمرین و تست کمتر از ۲ درصد است که می‌توان نتیجه گرفت مدل اوورفیت نشده و داده‌های تمرین را حفظ نکرده است و دارای تعمیم‌پذیری خیلی خوبی است. همچنین معیارهای ارزیابی تقریباً مثل هم هستند که این موضوع نشان می‌دهد مدل توانایی بالایی در تشخیص صحیح کلاس‌های مثبت و منفی دارد، بدون این‌که به نفع یک کلاس خاص متمایل شود.

Stacked Learners (۵)

Training ۱.(۵)

: GridSearch Parameters and Classifiers

```
classifiers_and_params = [  
    (RandomForestClassifier(), {'n_estimators': [50, 100], 'max_depth': [None, 5, 10],  
                                'min_samples_split': [5, 10], 'min_samples_leaf': [1, 2, 5]}),  
    (SVC(probability=True), {'C': [0.1, 1, 10], 'kernel': ['rbf'], 'gamma': ['auto', 'scale']}),  
    (LogisticRegression(), {'C': [0.1, 1, 10]}),  
    (KNeighborsClassifier(), {'n_neighbors': [3, 5, 7]})  
]
```

شکل ۲۰:

: Best Parameters

n_estimators: ۱۰۰، min_samples_split: ۵، min_samples_leaf: ۱، max_depth: ۱۰ : **RandomForest** •

kernel: rbf، gamma: auto، C: ۱۰ : **SVM** •

C: ۰.۱ : **LogisticRegression** •

n_neighbors: ۷ : **KNN** •

Model Evaluate ۲.(۵)

	Training	Test
Accuracy	0.895083	0.862000
Precision	0.895089	0.862060
Recall	0.895083	0.862000
F1 Score	0.895083	0.861999

شکل ۲۱:

با توجه به معیارهای ارزیابی میتوان دید کلسیفایر نهایی ما دارای دقت خوبی است و همچنین اختلاف بین معیارها در داده تست و تمرین حدود ۳ درصد است که میتوان تعمیم پذیری خوبی را نتیجه گرفت. با توجه به اینکه کلسیفایرهای بسیار متنوع هستند، برای دستیابی به بالاترین درصد ممکن میتوان کلسیفایرهای بیشتر برای Base Learners و کلسیفایرهای مختلف را به عنوان کلسیفایر نهایی امتحان کرد.

۳.۵ (Extra) overfitting

	Training	Test
Accuracy	0.912083	0.866667
Precision	0.912222	0.866945
Recall	0.912083	0.866667
F1 Score	0.912075	0.866652

	Training	Test
Accuracy	1.0	0.819333
Precision	1.0	0.819344
Recall	1.0	0.819333
F1 Score	1.0	0.819328

شکل ۲۳: بعد از جلوگیری از اورفیت

شکل ۲۲: قبل جلوگیری از اورفیت

در این قسمت یک درخت تصمیم اورفیت شده به کلسیفایر های قبلی اضافه کردیم. همانطور که در در عکس ها مشخص هست فقط با اووفیت شدن یک مدل، مدل نهایی اورفیت شده است و فاصله دقت بین داده آموزش و تست حدود ۱۹ درصد است. با اجرای الگوریتم برای جلوگیری از اورفیت فاصله دقت داده های آموزش و تست به ۵ درصد کاهش یافت و دقت روی داده های تست مثل قبل شد که نشان می دهد ما با موفقیت از اورفیت شدن مدل جلوگیری کرده ایم.