

به نام خدا
دانشگاه صنعتی امیرکبیر (پلی تکنیک تهران)
دانشکده مهندسی کامپیوتر



یادگیری ماشین

تکلیف چهارم

استاد درس: دکتر ناظر فرد

امیرحسین کاشانی

۴۰۰۱۳۱۰۷۱

amkkashani@gmail.com

نیم سال اول ۱۴۰۱-۱۴۰۲

Q1.....	3
a)	3
b)	3
c).....	3
d)	3
e)	4
f)	4
g)	5
h)	6
i)	7
j)	10
Q2.....	11
Q3.....	13
a)	13
b)	13
c).....	13
d)	13
Q4.....	14
a)	14
b)	14
c)	15
d).....	17
e)	19
f).....	20
g).....	21
h).....	21

Q1

a)

روش‌های ensemble روش‌هایی هستند که با ترکیب estimatorهای متخلف و رای گیری بین آن‌ها سعی در رسیدن به generalizability / robustness در مدل خود دارند. و به دو دسته boosting , averaging تقسیم می‌شوند.

در روش‌های Averaging که Bagging نیز عضوی از آن‌ها است مدل‌های متفاوت مستقلاً ساخته می‌شوند و در نهایت نظر آن‌ها باهم میانگین گرفته می‌شود. مدل نهایی عملکرد بهتری نسبت به هر یک مدل‌های ورودی خواهد داشت. (Averaging امکان موازی سازی را نیز دارد)

در روش‌های Boosting مدل به صورت sequential ایجاد می‌شود و امکان موازی سازی وجود ندارد و در هر گام سعی در بهبود نقاط یا فضایی از داده‌ها هستیم که در آن دقت یا کارایی مدل ضعیف بوده است. این روش عموماً از تجميع مدل‌های ساده بدست می‌آید.

**استفاده از Bagging در مواجهه با overfit کار آیی بیشتری نسبت به Boosting دارد و همچنین Bagging زمانی مناسب تر است که هر یک از base_estimatorها پیچیده به حساب بیایند ولی در Boosting عموماً base_estimatorها ساده تر و ضعیف تر هستند.

b)

در گام اول می‌توان گفت که این کار در زمانی که داده‌های کم باشند ایجاد مشکل می‌کند و با حذف داده‌ها در هر انتخاب داده‌های کمتری برای ایجاد مدل‌ها خواهیم داشت و عدم وجود اشتراک بین داده‌های مدل‌های مختلف باعث می‌شود که خرد جمعی رخ ندهد بلکه هریک از مدل‌ها می‌تواند بر روی یک قسمت نظر دهد.

در گام دوم که جواب اصلی به سوال می‌باشد این است که در صورتی که انتخاب بدون جایگذاری باشد احتمال ورود داده‌ها به مدل‌های مختلف دیگر توزیع uniform (یکنواخت) نخواهد داشت و شانس حضور در مدل‌ها متفاوت خواهد بود و توزیع داده‌ها در مدل ما به ترتیب ایجاد مدل‌ها وابسته خواهد شد.

c)

این کار چندان روش مناسبی نیست.

با استفاده از مدل‌های پیچیده تر در روش‌های Boosting در اصل داریم از ماهیت Ensemble دور می‌شویم زیرا در Boosting با تمرکز بر اشتباهات مدل سعی در کاهش bias داریم و استفاده از مدل پیچیده از ابتدا bias پایینی دارد و اشتباهاتی کمی از این رو یا تعداد iteration کمی را باید اتخاذ کنیم یا اینکه مدل به سرعت overfit خواهد کرد. از این رو در Boosting بهتر است که مدل‌ها ساده باشند.

d)

برای حل این مساله به شکل دقیق نیازمند این هستیم که توزیع احتمالاتی بر روی نظرات اشتباه را داشته باشیم از این رو فرض می‌کنیم که حالات خطا به شکل یکنواخت در estimatorهای ما پخش شده است یعنی احتمال اینکه لیبیل ۱ ، به اشتباه ۲ یا ۳ یا... تشخیص داده شود یکی است.

حال با این فرض اگر مدل ما صرفاً رندوم تصمیم گیری می کرد صحت آن برابر ۱۰۰ تقسیم بر ۷ می شد که یعنی ۱۴.۲ ، پس مدل های ما از رندوم نیز ضعیف تر هستند.

در این شرایط اگر با اغماض ۱۴.۲ را برابر ۱۴ فرض کنیم در بهترین حالت مدل های ما رندوم هستند و رای گیری بین یک سری مدل رندوم یک مدل رندوم دیگر ایجاد می کند که در بهترین حالت همان ۱۴ درصد صحت را بر می گرداند ولی اگر بخواهیم دقیق تر بررسی کنیم مدل های ما در جهت منفی دارند رای گیری می کنند یعنی همگرایی و رای گیری آن ها بدلیل اینکه احتمال پاسخ درست کمتر از میانگین احتمالات اختصاص داده است، یک پاسخ غلط را به عنوان پاسخ نهایی پر رنگ می کند از این رو در این مورد به هیچ وجه نباید از رای گیری استفاده نمود.

e)

با توجه به اینکه Boosting بصورت sequential در حال اجرا است و در هر iteration نیازمند آن است که بر روی تمامی داده ها حرکت کند پس گزینه مناسبی برای دیتاست های بزرگ نیست اما از طرف دیگر Bagging می تواند بصورت موازی آموزش یابد و در نتیجه از لحاظ زمانی بهتر عمل می کند.

برای ابعاد زیاد نیز می توان گفت که مدل های ساده در نظر گرفته شده در Boosting شاید به نظر بیاید که ضعیف عمل کنند اما با توجه به اینکه می توان تعداد iteration ها را بیشتر کرد بر این مساله نیز غلبه می کنند. در نتیجه خواهیم داشت

دیتاست بزرگ و ابعاد کم \leq Bagging

دیتاس کوچک و ابعاد بزرگ \leq Boosting

f)

نظری ندارم

g)

g)

برای حل این مسأله باید معادله‌ی زیر برقرار باشد.

$$\begin{cases} \nabla f(x, y) = \lambda \nabla g(x, y) \\ g(x, y) = c \end{cases}$$

g-1)

$$\begin{cases} -2x = \lambda \times 1 \\ -4y = \lambda \times 1 \end{cases} \rightarrow \begin{cases} x = -\lambda/2 \\ y = -\lambda/4 \end{cases} \quad \textcircled{I} \quad \begin{matrix} x = 1/2 \\ y = 1/4 \end{matrix}$$

$$x + y = 1 \xrightarrow{\textcircled{I}} -\frac{1}{2}\lambda = 1 \rightarrow \lambda = -\frac{2}{1} \quad \textcircled{I}$$

$$f(1/2, 1/4) = 2 - 1/4 - 1/4 = 1/2$$

g-2)

$$\begin{cases} 2x = \lambda \times 1 \\ 2y = \lambda \times 1 \end{cases} \rightarrow \begin{matrix} x = y \\ x = -1 \\ y = -1 \end{matrix} \quad f(-1, -1) = 2$$

g-۳)

$$\left. \begin{array}{l} r_x^2 = r_x \lambda \\ r_y^2 = r_y \lambda \end{array} \right\} \rightarrow \begin{array}{l} x=0 \\ y=0 \end{array} \rightarrow \begin{array}{l} x=y \rightarrow r_x = r_y \rightarrow x = \frac{r_x}{3} \\ y = \frac{r_y}{3} \end{array}$$

$$x^2 - 1 \geq 0 \rightarrow \frac{4}{9} \lambda^2 \geq 1 \rightarrow \begin{cases} \lambda > \frac{3}{2} \\ \lambda < -\frac{3}{2} \end{cases}$$

محدودی نتوانند $-\infty$ و $+\infty$ باشند و معادلات پاسخ خاصی به

$$f(-\infty, -\infty) = -\infty$$

ساختی دهند.

h)

h)

با توجه به جابجایی نقاط

معادله خط به دست آمده از SVM برابر زیر می باشد

$$x_2 = -x_1 + 3$$

زیرا باید بر $x_2 - x_1$ عمود بوده و از میانگین آنها نیز عبور کند.

i)

i)

$$\text{kernel matrix} \rightarrow \begin{cases} K = K^T & \text{متقارن} \\ \alpha^T K \alpha \geq 0 & \text{Positive - semi - definite} \end{cases}$$

$$i-1) K(x, z) = K_1(x, z) + K_2(x, z)$$

جمع دو ماتریس متقارن، متقارن خواهد شد و اگر هر دو مثبت نیمه تعریف باشند، مثبت نیمه تعریف خواهند بود.
برای اثبات این موضوع.

CS Scanned with CamScanner

$$z^T K z = z^T (K_1 + K_2) z = z^T K_1 z + z^T K_2 z \quad (I)$$

$$z^T K_1 z \geq 0 \quad \left. \begin{array}{l} z^T K_2 z \geq 0 \end{array} \right\} \rightarrow z^T K z \geq 0 \quad \checkmark \text{ positive semi-definite}$$

$$z^T K_2 z \geq 0$$

(I)

∴ کرنل K مثبت است

$$i-2) K(x, z) = K_1(x, z) - K_2(x, z)$$

تقارن مشابه مثال قبل برقرار است. و K_1 positive semi-definite خواهد بود

مثال تقارن K کرنل نیست

$$K_1 \quad K_2 \quad K$$

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} - \begin{bmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 2 \end{bmatrix} = \begin{bmatrix} -1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & -1 \end{bmatrix}$$

$$i-3) K(x, z) = \alpha K_1(x, z)$$

شرط تقارن با توجه به اینکه α اسکالری باشد قطعی شود

$$z^T K z = z^T (\alpha K_1) z = \alpha z^T K_1 z \rightarrow \alpha z^T K_1 z \geq 0$$

$$z^T K_1 z \geq 0 \quad \left| \begin{array}{l} \text{که در صورتی که } \alpha \geq 0 \\ \text{کرنل معتبر است} \end{array} \right.$$

PAPCO

$$\alpha = -1$$

مثال نقض برای حالتی که $\alpha < 0$ باشد

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \times \alpha = \begin{bmatrix} -1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & -1 \end{bmatrix}$$

K_1

K

$$i - 4) \quad K(x, z) = K_1(x, z) K_2(x, z)$$

راه حل زیر مرتبه‌ی راه حل یافت شده است می باشد

$$\text{بوجه فرض} \quad \begin{cases} K_1(x, z) = (\phi_{(x)}^{(1)})^T (\phi_{(z)}^{(1)}) \\ K_2(x, z) = (\phi_{(x)}^{(2)})^T (\phi_{(z)}^{(2)}) \end{cases}$$

$$K(x, z) = K_1(x, z) K_2(x, z)$$

$$= \sum_i \phi_i^{(1)}(x) \phi_i^{(1)}(z) \sum_j \phi_j^{(2)}(x) \phi_j^{(2)}(z)$$

$$= \sum_i \sum_j \phi_i^{(1)}(x) \phi_i^{(1)}(z) \phi_j^{(2)}(x) \phi_j^{(2)}(z)$$

$$= \sum_i \sum_j [\phi_i^{(1)}(x) \phi_j^{(2)}(x)] [\phi_i^{(1)}(z) \phi_j^{(2)}(z)] \quad \text{def } \psi_{(i,j)} = \phi_i^{(1)}(x) \phi_j^{(2)}(x)$$

$$= \sum_{(i,j)} \psi_{(i,j)}(x) \psi_{(i,j)}(z) \rightarrow K(x, z) = \psi(x)^T \psi(z)$$

K گزین است

j)

در حالت عادی یک داده زمانی منجر به محاسبه مجدد مرزها می‌شود که بین داده‌های مرزی بیفتد به بیان دیگر در جاده ی بین داده‌های دو کلاس قرار گیرد. برای این حالت می‌توانیم یک شرط کنترل قرار دهیم که در صورتی که میزان خطای سیستم از یک threshold کمتر بود محاسبه مجدد صورت نگیرد و در مواردی که اشتباه کم است اغماض کند .

K)

این شروط به ما اجازه می‌دهد تا در حل مسائل برنامه نویسی غیر خطی بتوانیم مانند روش لاگرانژ از شروطی به جز شرط تساوی نیز استفاده کنیم. هنگامی که مسئله اولیه محدب باشند شرایط KKT برای نقاط بهینه مسئله اولیه و مسئله دوگان صادق هستند، یا به عبارت دیگر فاصله دوگانی صفر می‌باشد.

تعریف مساله

$$\begin{aligned} & \text{minimize} && f(x) \\ \text{mtfc subject to} && g_i(x) \leq 0; & i \in \{1, \dots, m\} \\ && h_j(x) = 0; & j \in \{1, \dots, \ell\} \end{aligned}$$

شروط KKT

۱. مساله دارای جواب باشد

۲. مساله دارای دوگان باشد

۳. شرط Complementary slackness

$$\mu_i g_i(x^*) = 0, \text{ for } i = 1, \dots, m.$$

۴. شرط ایستا برقرار باشد

$$-\nabla f(x^*) = \sum_{i=1}^m \mu_i \nabla g_i(x^*) + \sum_{j=1}^{\ell} \lambda_j \nabla h_j(x^*),$$

با در نظر گرفتن شروط فوق و حل معادلات بدست آمده می‌توانیم مینم تابع مورد نظر را بدست آوریم.

Q2

نتایج بدست آمده نشان می‌دهد که همه کرنل‌های ساخته شده در مجموعه داده آپریس عملکرد مشابهی داشته‌اند البته این تشابه را می‌توان به وجود کم داده‌های تست نیز مرتبط دانست. اما کرنل sigmoid در حالتی که به تنهایی گرفته شود عملکرد به شدت ضعیفی از خود نشان داده است.

```
kernel 1
```

```
accuracy : 0.9666666666666667
```

```
confusion matrix
```

```
[[10  0  0]
```

```
 [ 0 13  0]
```

```
 [ 0  1  6]]
```

```
kernel 2
```

```
accuracy : 0.9666666666666667
```

```
confusion matrix
```

```
[[10  0  0]
```

```
 [ 0 13  0]
```

```
 [ 0  1  6]]
```

```
kernel 3
```

```
accuracy : 0.9666666666666667
```

```
confusion matrix
```

```
[[10  0  0]
```

```
 [ 0 13  0]
```

```
 [ 0  1  6]]
```

kernel 4

accuracy : 0.9666666666666667

confusion matrix

```
[[ 9  0  0]
 [ 0  9  0]
 [ 0  1 11]]
```

RBF only

accuracy : 0.9666666666666667

confusion matrix

```
[[10  0  0]
 [ 0 13  0]
 [ 0  1  6]]
```

sigmoid only

accuracy : 0.23333333333333334

confusion matrix

```
[[ 0  0 10]
 [ 0  0 13]
 [ 0  0  7]]
```

Q3

a)

در این بخش random_state در بخش split و make_classification برابر ۴۲ قرار داده شده است.

```
1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
2 X_train.shape, X_test.shape, y_train.shape, y_test.shape

((80000, 2), (20000, 2), (80000,), (20000,))
```

b)

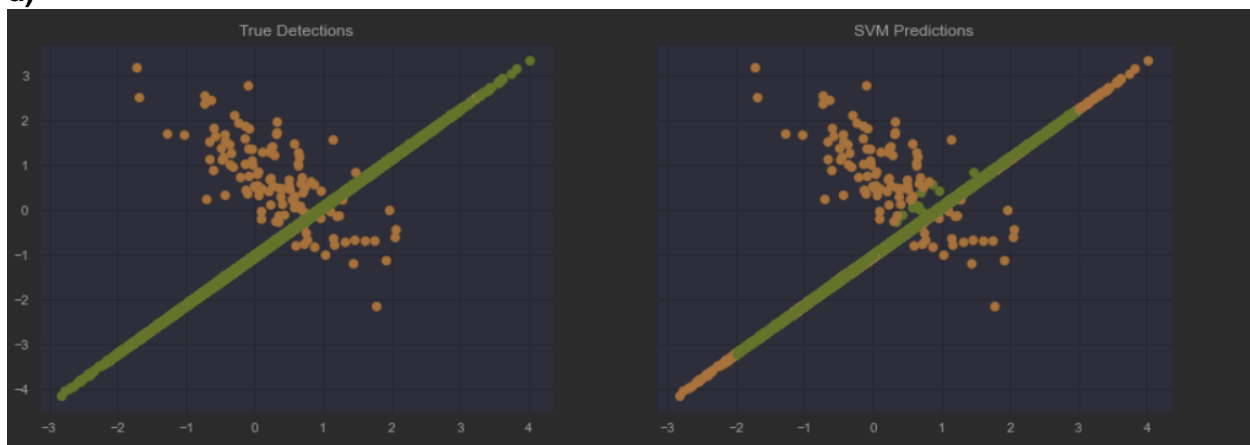
One_class_svm یک روش Unsupervised مبتنی بر SVM است که با به کارگیری SVM به دنبال ایجاد یک ابر صفحه برای تشخیص داده‌های پرت می‌باشد ملاک این دسته بندی تراکم داده‌ها و یادگیری توزیع آماری داده‌ها می‌باشد و بر اساس آن به فضاوت در رابطه با اینکه آیا این داده یک داده ی پرت است یا خیر می‌پردازد.

c)

با توجه به اینکه تعداد داده‌های پرت ما دسته به شدت کوچکی از داده‌های ما است پس هدف ما شناسایی هر چه بیشتر آن‌ها می‌باشد از این رو معیاری مانند Recall می‌تواند گزینه ای مناسب برای ارزیابی الگوریتم ما باشد.

	precision	recall	f1-score	support
0	0.99	0.98	0.99	1976
1	0.06	0.08	0.07	24
accuracy			0.97	2000
macro avg	0.52	0.53	0.53	2000
weighted avg	0.98	0.97	0.98	2000

d)



Q4

a)

فرض بر این است که سمت مثبت، سالم بودن نرم افزار در نظر گرفته شده است.

از دید امنیتی قطعاً باید **False positive rate** کاهش یابد (مهم ترین معیار است) زیرا این که بتوانیم همه بدافزارها را شناسایی کنیم و به اشتباه تعدادی بد افزار را به عنوان نرم افزار سالم پیش بینی نکنیم اولویت اول است هر چند که این مهم در صورتی که بدست بیاید که تعدادی نرم افزار سالم را نیز با سخت گیری بیش از حد مخرب پیش بینی کنیم که False Negative را افزایش می دهد.

اما در سمت user باید هر دو مولفه در کنار یک دیگر وجود مهم هستند اما False Negative Rate اهمیت بیشتری در تجربه کاربری دارد زیرا در صورتی که False Negative افزایش یابد سیستم ما به اشتباه تعدادی از برنامه هایی که کاربر به آن ها نیاز داشته است را مخرب شناسایی می کند مثل windows defender که همیشه فقط فایل های کرک را به عنوان مخرب شناسایی می کند از این رو False Negative Rate در صورتی که از حد استاندارد بیشتر باشد برای تجربه کاربری لطمه بزرگ تری به حساب می آید.

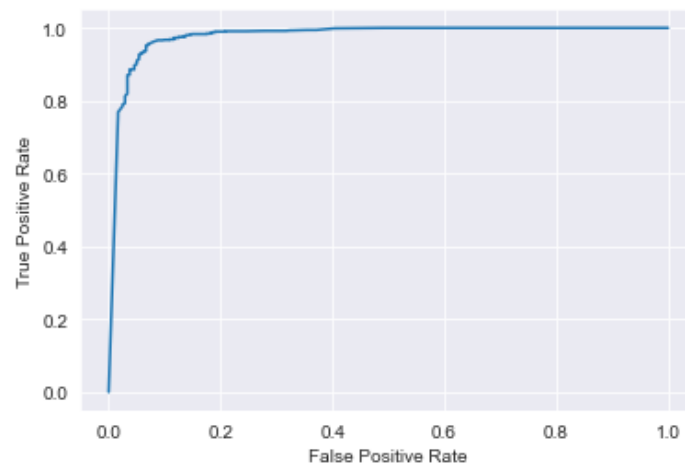
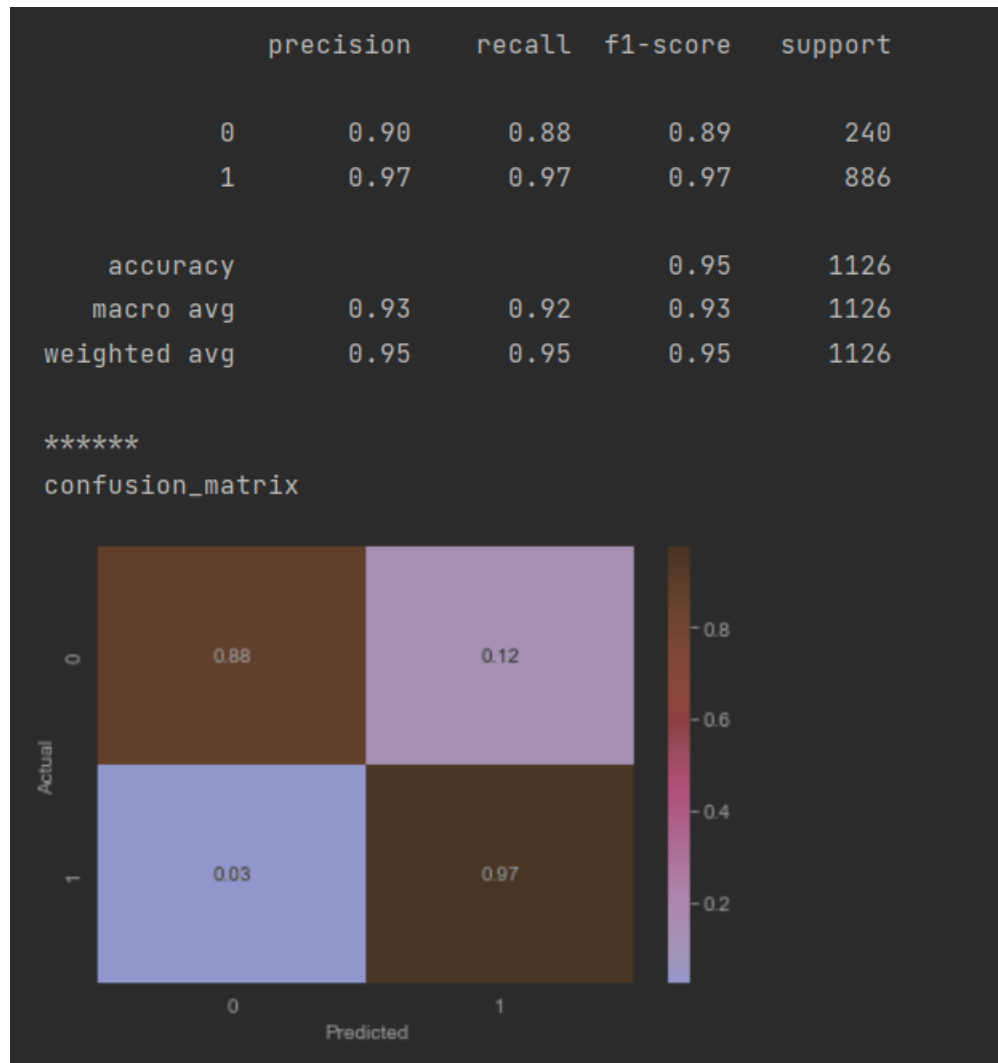
b)

stratified sampling

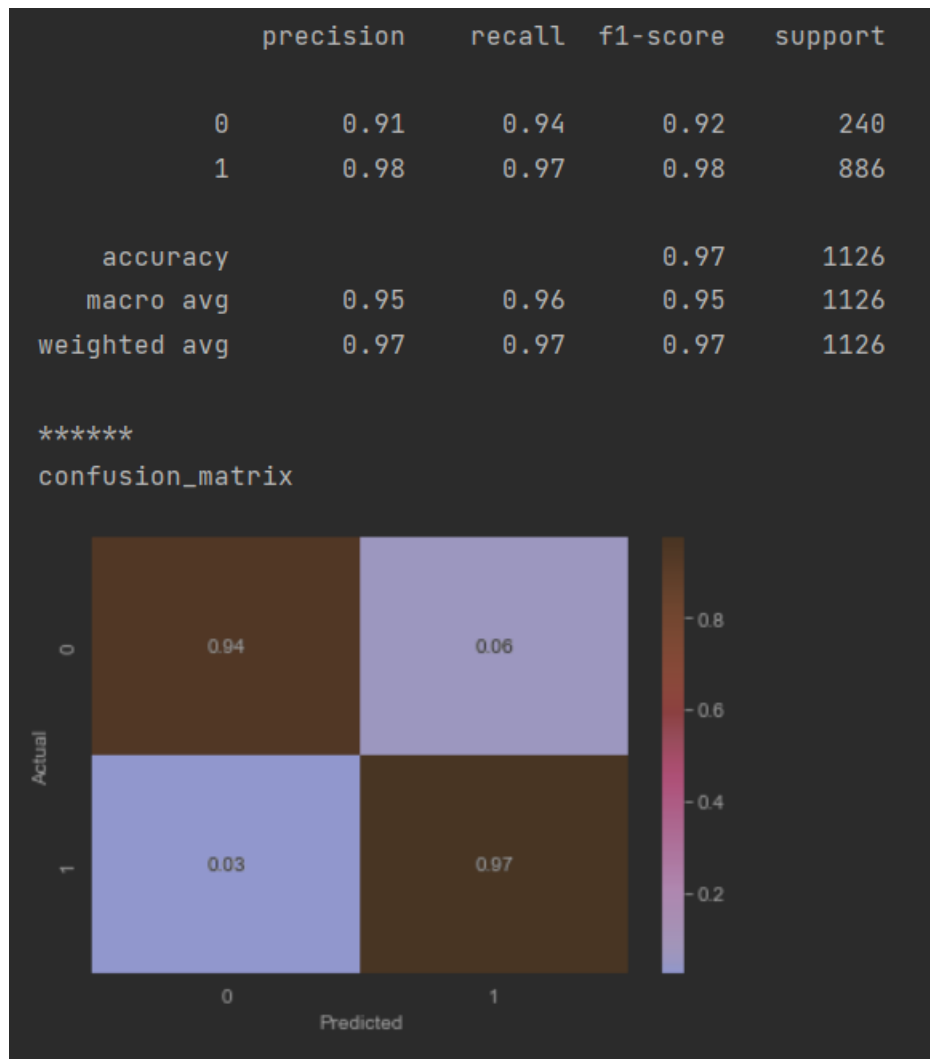
روشی است که در آن به جای نمونه برداری رندوم از داده ها به شکلی نمونه برداری می کنیم که نسبت توزیع داده ها در هر یک از دسته ها ثابت بماند به این معنی که اگر ۱ درصد از داده ها در داده های موجود مربوط به کلاس C هستند بعد از جدا کردن داده test, train, در هر کدام یک درصد از داده ها مربوط به گروه C می باشد. کاربرد این روش بیشتر در داده هایی است که تعداد داده های یک دسته یا چند دسته به کم باشد یا به طور کلی دیتاست کوچکی داشته باشیم و در نتیجه در صورت عدم استفاده از این روش ممکن است جمعیت برخی از دسته ها به صفر میل کنند.

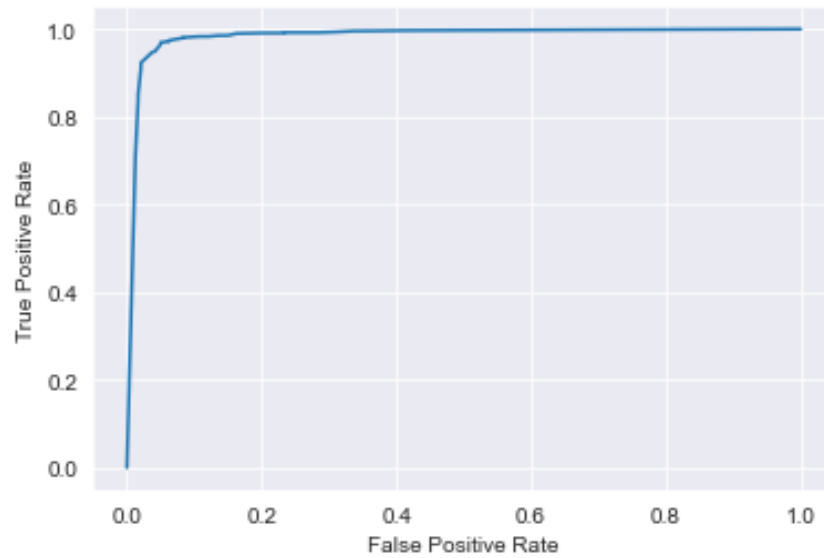
برای جدا کردن داده ها به داده های آموزش و تست فیلد stratified در تابع train_test_split استفاده شده است. در ادامه در بخش e از تابع StratifiedShuffleSplit برای این کار استفاده گردیده است.

c)
bagging with KNN



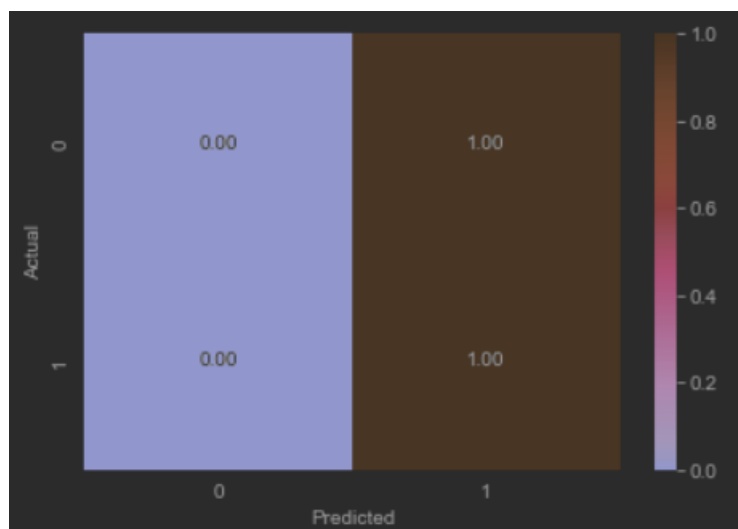
Bagging with tree

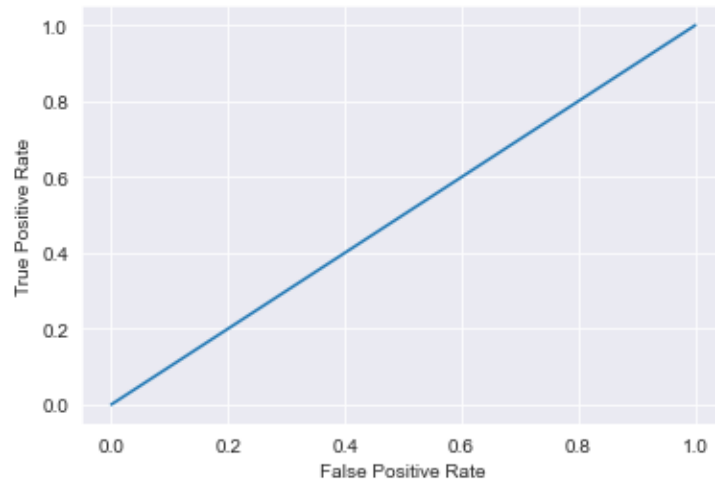




d)
Adaboost SVC

	precision	recall	f1-score	support
0	0.00	0.00	0.00	240
1	0.79	1.00	0.88	886
accuracy			0.79	1126
macro avg	0.39	0.50	0.44	1126
weighted avg	0.62	0.79	0.69	1126

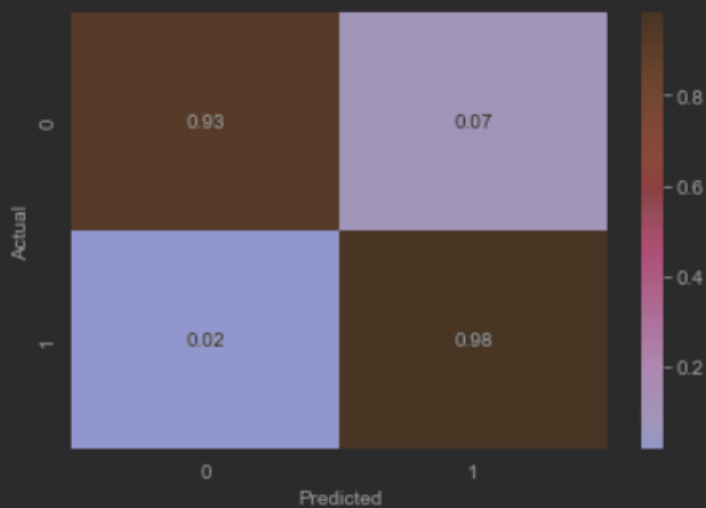


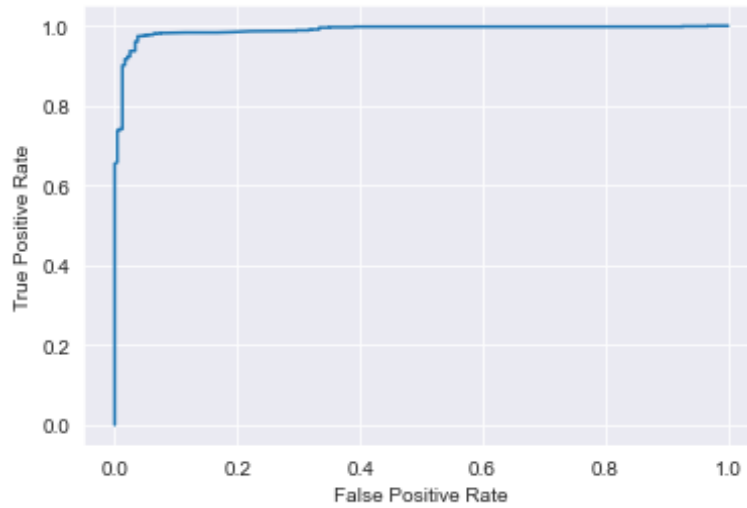


Ada Boost tree

	precision	recall	f1-score	support
0	0.93	0.93	0.93	240
1	0.98	0.98	0.98	886
accuracy			0.97	1126
macro avg	0.96	0.95	0.95	1126
weighted avg	0.97	0.97	0.97	1126

confusion_matrix

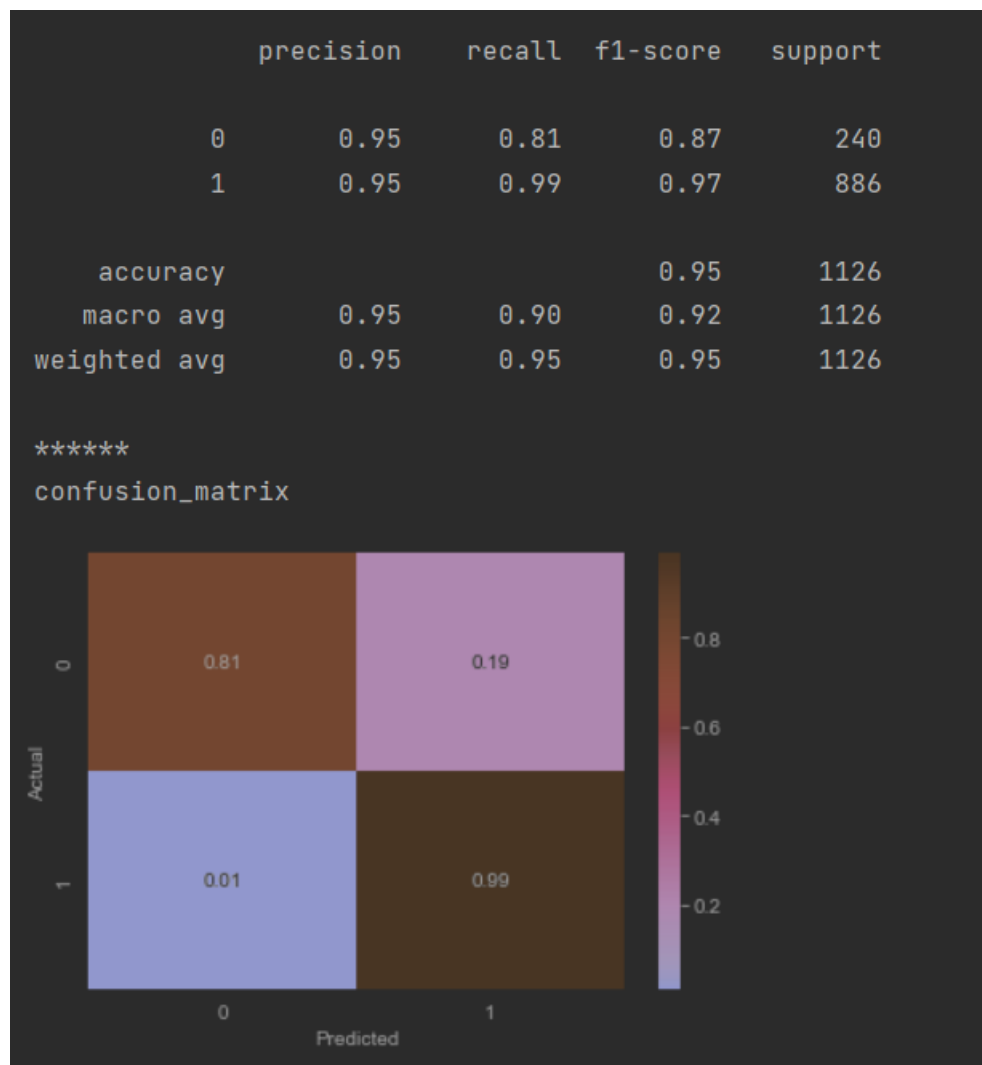




e)

در این بخش از تابع `StratifiedShuffleSplit` برای ایجاد نمونه‌ها استفاده کرده ایم و ۱۰۰ بار این تابع را فراخوانی کرده و در هر بار ۵ درصد داده را برداشته ایم.

f)



g)

یک نمونه از ماتریس ۱۰۱ بعدی ایجاد شده

```
array([[1, 1, 1, ..., 1, 1, 1],
       [1, 1, 1, ..., 1, 1, 1],
       [1, 1, 1, ..., 1, 1, 1],
       ...,
       [0, 1, 0, ..., 1, 0, 0],
       [1, 1, 1, ..., 1, 1, 1],
       [1, 1, 1, ..., 1, 1, 1]], dtype=int64)
```

h)

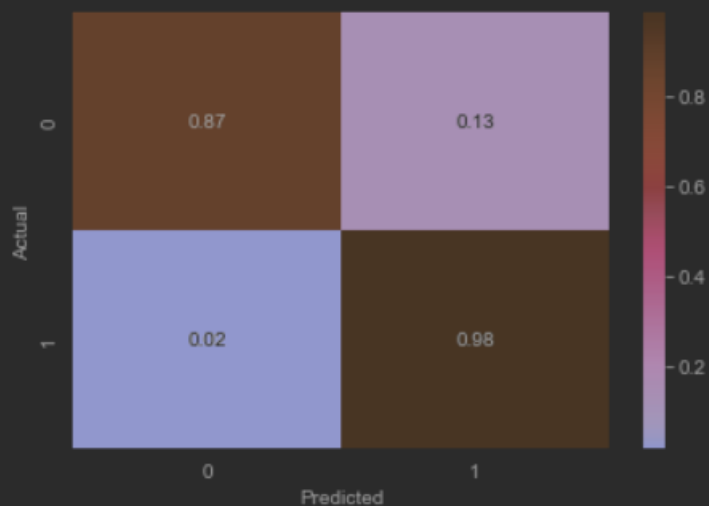
نتیجه svc

```
              precision    recall  f1-score   support

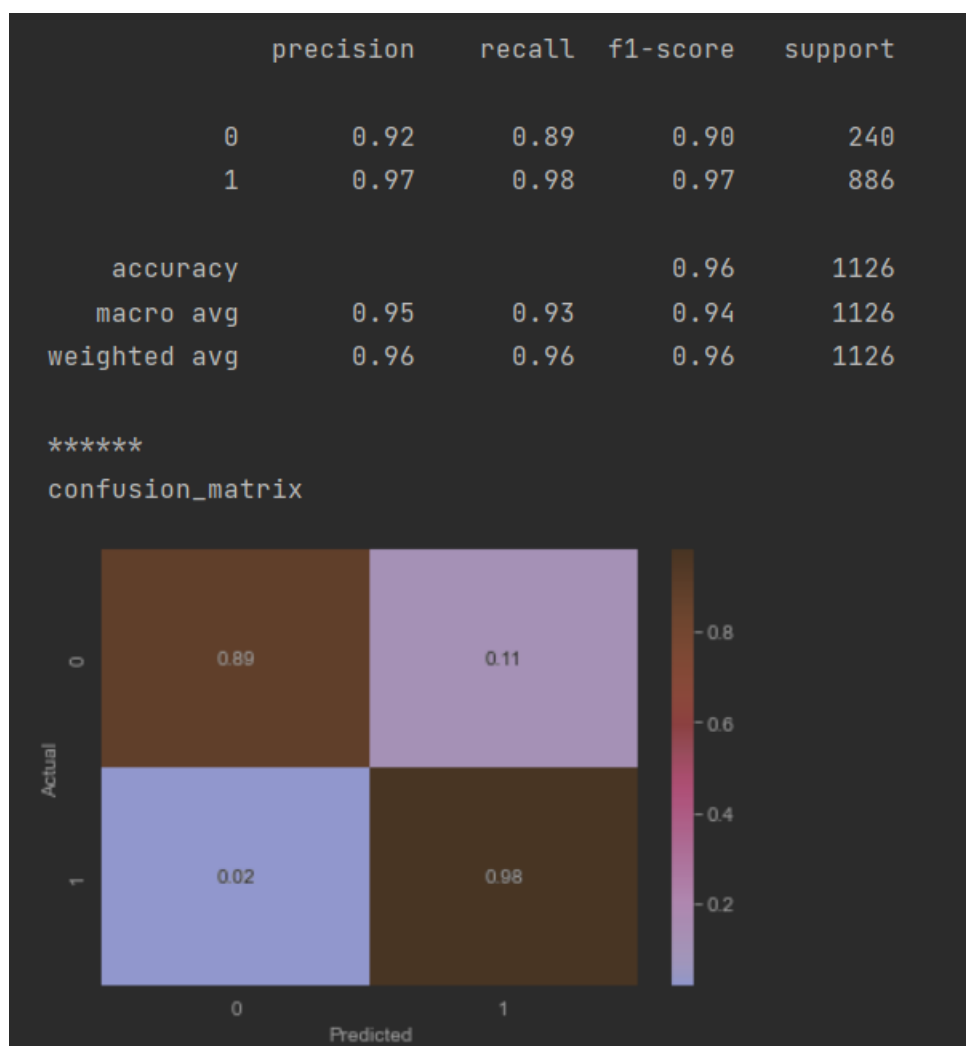
     0       0.93      0.87      0.90        240
     1       0.97      0.98      0.97        886

 accuracy              0.96        1126
 macro avg           0.95      0.93      0.94        1126
 weighted avg        0.96      0.96      0.96        1126
```

confusion_matrix



نتیجه decision tree



به طور کلی می توان گفت که نتیجه بدست آمده در همه مدل ها قابل قبول است و مقایسه روش ها زمانی که همگی در بازه ۹۵ به بالا هستند زیاد عادلانه نیست و دلیل برتری آن ها می تواند صرفا به خاص بودن اون دیتاسیت مربوط باشد. اما به طور کلی مدل هایی که در بخش آخر تولید شده بودند توزیع اعداد در confusion ماتریس در مجموع بهتر بود به این معنی که FN یا FP وجود ندارد که خیلی متفاوت باشد. در بخش آخر نیز عملکرد درخت تصمیم بهتر از SVC بود.