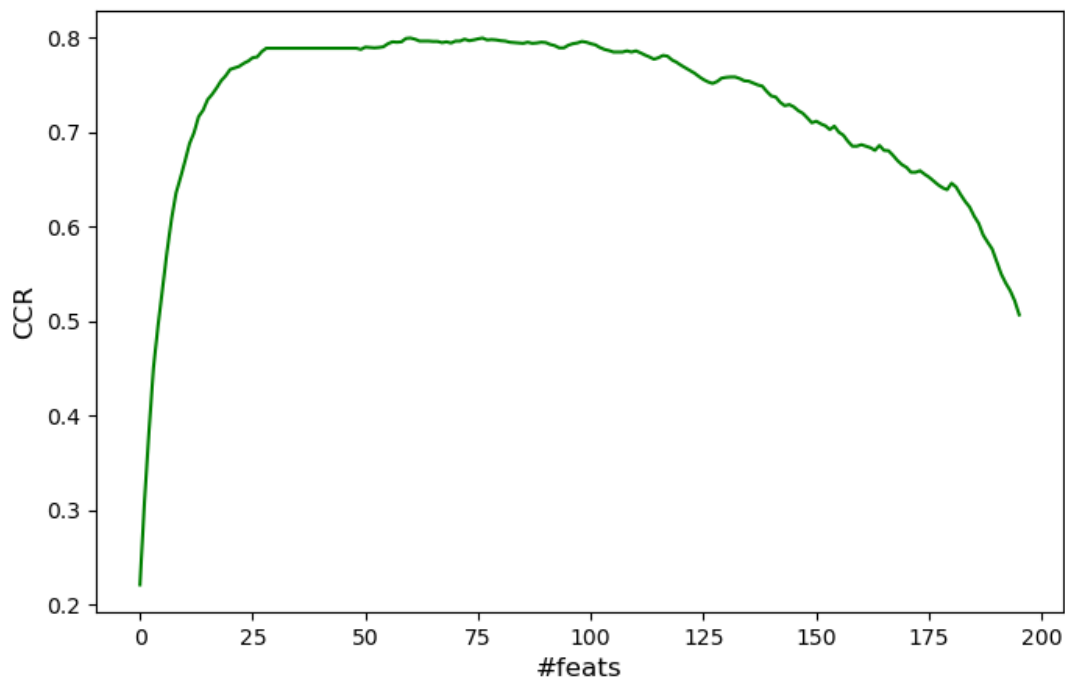


۱- من این دو الگوریتم رو بر روی کل دیتاست (یعنی هر ۱۹۶ feature) اجرا کردم. نتایج حاصل از اجرای هر یک به صورت زیر می‌باشد:

نتیجه‌ی Forward Selection:



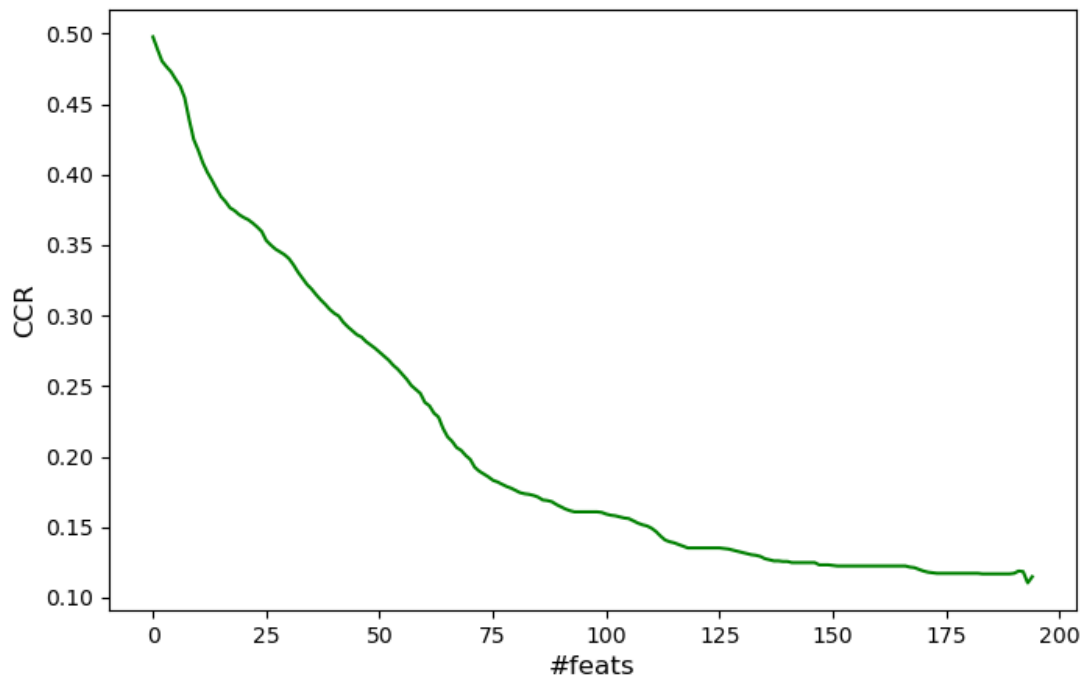
مقادیر CCR بیشینه و تعداد featureهای انتخابی که به ازای آن maximum CCR حاصل می‌شود به صورت زیر می‌باشد:

0.7996 61

یعنی میزان maximum CCR با استفاده از Forward Selection برابر 0.7996 بود و تعداد ویژگی‌هایی که با این الگوریتم منجر به این maximum CCR می‌شود برابر است با 61 ویژگی. و شماره‌ی ویژگی‌های انتخاب شده به صورت زیر می‌باشد:

```
[35, 102, 91, 106, 131, 89, 49, 104, 47, 65, 132, 63, 74, 133, 52, 162, 135, 76, 77, 78, 144, 149, 116, 61, 118, 64, 90, 75, 148, 0, 1, 2, 3, 12, 13, 14, 15, 27, 28, 84, 98, 154, 167, 168, 181, 182, 183, 194, 195, 120, 51, 134, 147, 80, 101, 103, 46, 60, 79, 48, 146]
```

نتیجه‌ی Backward Selection:

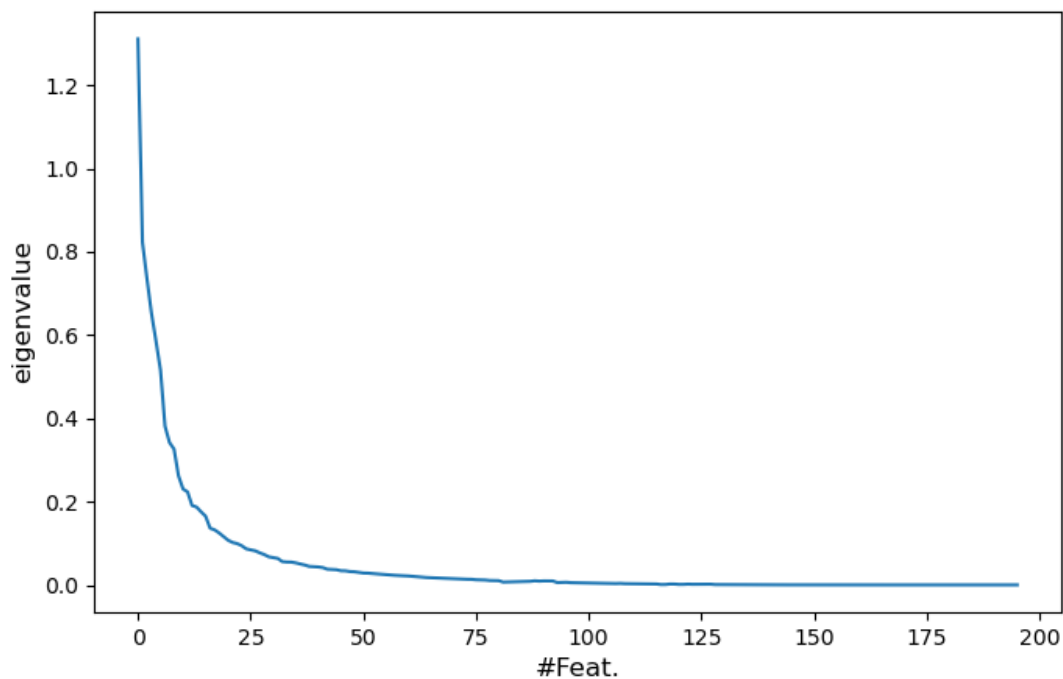


0.4976 196

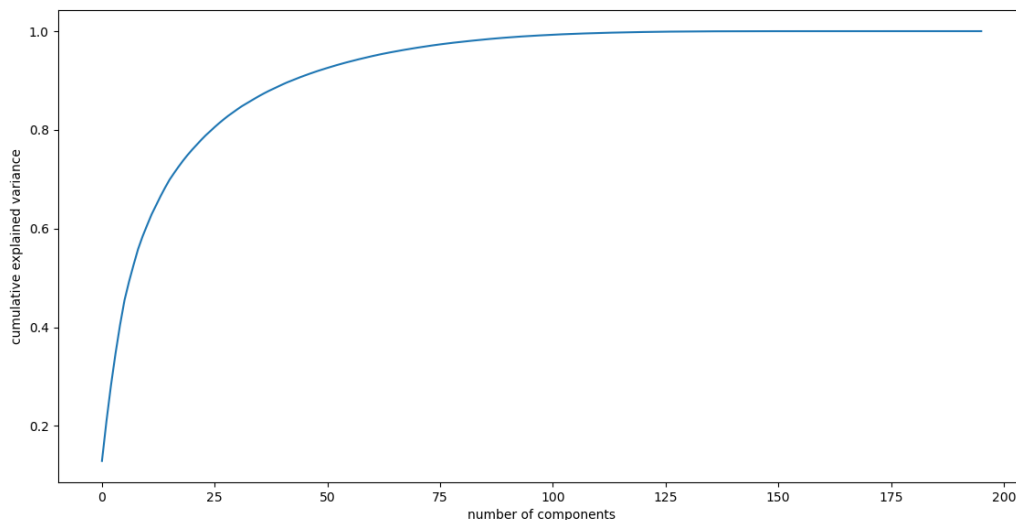
بیشترین میزان CCR برابر با 0.4976 می باشد که با انتخاب تمام ویژگی ها حاصل می شود. علت این تفاوت بارز در نتیجه ی FS و BS را می توان این گونه توضیح داد که: در backward selection مشکلی مشابه با مشکلی که در FS وجود داشت داریم. در BS پس از حذف ویژگی ها در مراحل اولیه در مراحل بعدی نمی توانیم ویژگی های حذف شده را reevaluate کنیم. و فقط می توانیم از مجموعه ویژگی های فعلی ویژگی کم کنیم.

۲- داده های train و test را split می کنیم. حال می خواهیم ماتریس کوواریانس را برای داده های train محاسبه کنیم. برای اینکار لازم است که داده ها centered شوند. یعنی میانگین هر ستون از ماتریس داده ها را محاسبه کرده (یعنی مقدار میانگین داده ها مربوط به هر feature) و سپس این ماتریس میانگین را از هر داده (هر سطر) کم می کنیم. با استفاده از تابع cov از پکیج numpy ماتریس کوواریانس را حساب کرده و با استفاده از تابع eig مقادیر و بردارهای ویژه ی این ماتریس را بدست می آوریم.

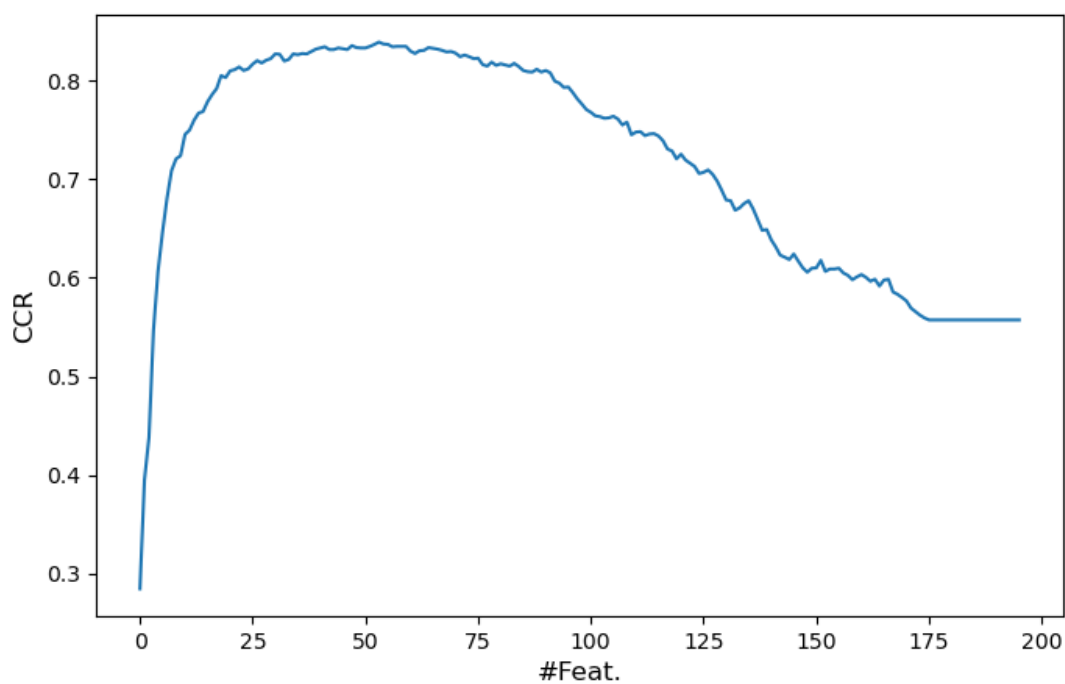
نمودار مقادیر ویژه ماتریس کواریانس بر حسب تعداد ویژگی به صورت زیر خواهد بود:



اما برای اینکه متوجه شویم که با چه تعداد ویژگی بیشترین مقدار ویژه حاصل می‌شود (در PCA به دنبال یافتن بردار ویژه‌ای هستیم که منتظر با بیشترین مقدار ویژه می‌باشد)، باید نمودار explained variability تجمعی بر حسب تعداد ویژگی را رسم کنیم. این نمودار به ما می‌گوید که به ازای چه تعدادی از ویژگی‌ها بیشتر واریانس موجود در داده‌ها قابل توضیح است. این نمودار به صورت زیر خواهد بود:



به طور معمول درصد explained variance در حدود 80% کافی و مطلوب است. این میزان واریانس تجمعی با تعداد ویژگی حدود ۳۰ حاصل می شود. اما اگر بخواهیم این موضوع را با نمودار CCR بر حسب تعداد ویژگی بررسی کنیم، خواهیم داشت:



می بینیم که نتیجه گیری بالا با این نمودار هم همخوانی دارد. ما در حدود ۳۰ ویژگی، به مقدار حدود 0.8 برای CCR می رسیم. البته می بینیم که با تعداد حدوداً نزدیک 60 ویژگی بیشترین مقدار CCR را خواهیم داشت.

این نتایج به صورت شهودی از روی نمودارها بود. حال با استفاده از کد مقدار بیشینه CCR و تعداد ویژگی‌های بهینه را بدست می‌آوریم:

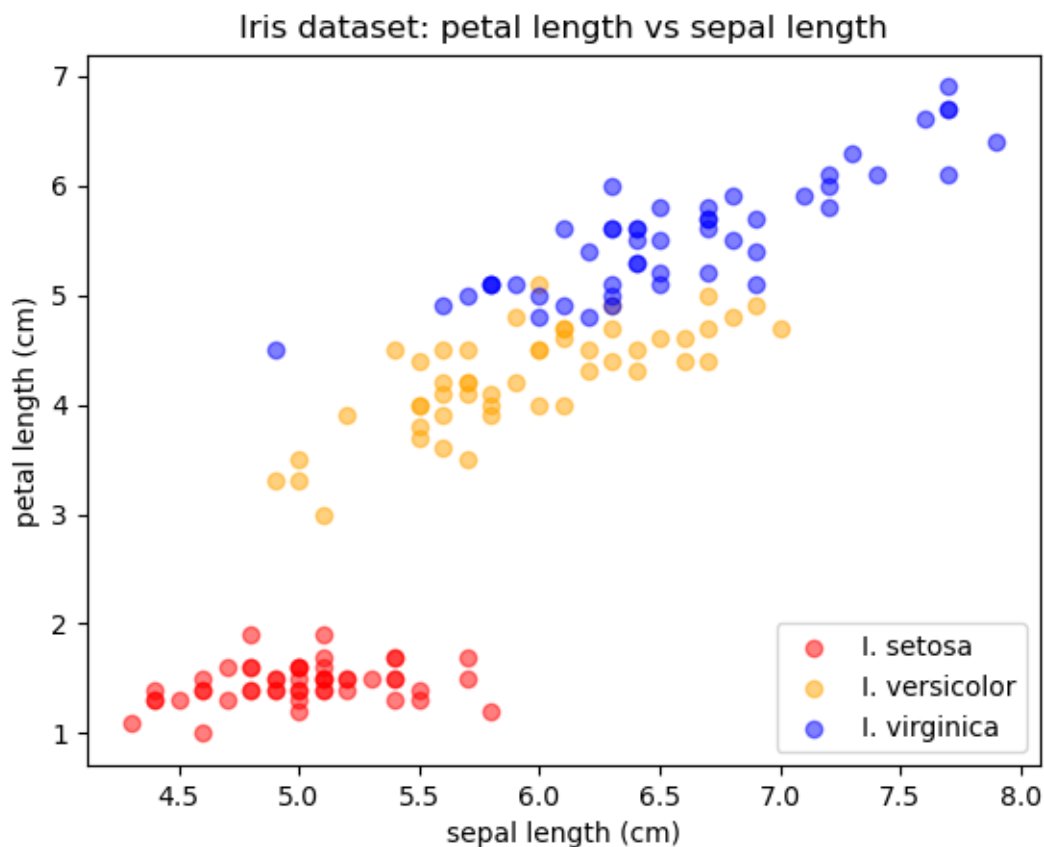
```
maximum CCR: 0.8388
optimum number of features: 53
```

می‌بینیم که نتیجه‌ی بدست آمده برای تعداد بهینه ویژگی‌ها و بیشترین مقدار CCR با تحلیل شهودی ما همخوانی دارد. (تقریباً) مقایسه با سوال ۱:

تعداد ویژگی‌های بدست آمده با الگوریتم Forward Selection (FS) در سوال ۱ برابر با 61 بوده که این تعداد ویژگی منجر به $CCR = 0.7996$ می‌شد. می‌بینیم که مقادیر به دست آمده برای CCR و تعداد بهینه‌ی ویژگی در این بخش اختلاف کمی با نتایج بدست آمده از FS دارند. دلیل این اختلاف می‌تواند این باشد که در الگوریتم FS ما امکان حذف ویژگی‌های انتخاب شده در مراحل قبل را نداریم. به عبارتی شاید مجموعه ویژگی‌های بهینه اصلاً شامل تک ویژگی‌ای که منجر به بیشترین CCR نسبت به تک تک ویژگی‌های دیگر می‌شود نباشد. اما در FS ما حتماً این ویژگی را در مرحله اول انتخاب می‌کنیم پس با اینکار امکان رسیدن به maximum CCR را به ازای مجموعه ویژگی‌های انتخابی خود می‌گیریم.

-۸

Scatter plot داده‌های کلاس‌های مختلف بر حسب دو ویژگی petal length و petal width به صورت زیر می‌باشد:



الف) حال SVM را با kernel های rbf و linear و polynomial اعمال می کنیم و با استفاده از 10-fold cross validation نتایج آنها را نشان می دهیم (مقدار $f1_score$ و ماتریس آشفستگی). در خروجی یک ماتریس سطری با اندازه ۳ داریم که درایه های آن نشان دهنده $f1_score$ برای کلاس های ۰ و ۱ و ۲ می باشد. (به ترتیب کلاس های setosa و versicolor و virginica). و یک ماتریس ۳ در ۳ داریم که همان ماتریس آشفستگی می باشد.

برای linear kernel داریم:

```
linear
[1.  0.94 0.94]
[[50  0  0]
 [ 0 47  3]
 [ 0  3 47]]
```

برای polynomial kernel داریم:

```
poly
[1.  0.92 0.92]
[[50  0  0]
 [ 0 46  4]
 [ 0  4 46]]
```

برای rbf kernel داریم:

```
rbf
[1. 0.90196078 0.89795918]
[[50  0  0]
 [ 0 46  4]
 [ 0  6 44]]
```

در مورد این kernel می توان گفت:

تفاوت این kernel ها در ساخت decision boundary بین کلاس ها می باشد. به طور کلی کاربرد kernel function این است که دیتاست اصلی را به فضایی با تعداد ابعاد بیشتر ببرند به این امید که داده ها در آن فضای با بعد بیشتر خطی شوند. کاری که RBF kernel انجام می دهد این است که از یک ترکیب غیرخطی از ویژگی ها استفاده می کند تا داده ها را به یک فضای با ابعاد بزرگتر برده که بتوان داده ها را در آن فضا با linear decision boundary از هم جدا کرد. پس می توان نتیجه گرفت که linear kernel و polynomial kernel از نظر پیچیدگی زمانی نسبت به RBF به صرفه ترند اما دقت کمتری نیز دارند.

اما به طور کلی هیچ گارانتی ای وجود ندارد که یک kernel عملکرد بهتری از دیگری دارد. و باید بر اساس مسئله و دیتاست آن تصمیم بگیریم که کدام کرنل بهتر است و از آن استفاده کنیم.

برای دیتاست iris در این تمرین می بینیم که linear kernel عملکرد بهتری دارد.

(ب) توضیح راجب پارامتر gamma و C:

پارامتر گاما مشخص می کند که اثر فاصله یک داده ی train تا چه حدی می تواند برسد. در واقع مقدار پارامتر گاما مشخص کننده ی میزان انحنایی است که می خواهیم decision boundary ما داشته باشد. مقادیر بالای گاما منجر به مدل های با بایاس بالا و واریانس کم (یعنی decision boundary با انحنای زیاد) و مقادیر پایین گاما منجر به مدل های با بایاس کم و واریانس زیاد (یعنی decision boundary با انحنای کم) می شود.

پارامتر C یک جریمه به ازای هر misclassified data point اضافه می کند. اگر مقدار C کم باشد، جریمه به ازای نقاط به اشتباه طبقه بندی شده نیز کم خواهد بود در نتیجه یک مرز تصمیم با مقدار margin بزرگتر انتخاب می شود اما در ازای آن تعداد نقاط به اشتباه طبقه بندی شده (misclassified points) بیشتر خواهد شد. اما اگر مقدار C بزرگ باشد، SVM سعی می کند تعداد نقاط به اشتباه طبقه بندی شده را کمینه کند. چرا که به ازای مقادیر بالای C، جریمه به

ازای نقاط به اشتباه طبقه‌بندی شده زیاد خواهد بود و در نتیجه decision boundary حاصل مقدار margin کوچکتری دارد.

:Gamma

خروجی(ماتریس سطری از f1_score ها که در سوال ۲ هم توضیح داده شد و ماتریس آشفستگی) به ازای مقادیر مختلف گاما. در خروجی ابتدا مقدار گاما چاپ شده و سپس f1_score کلاس‌ها و بعد هم ماتریس آشفستگی:

```
0.01
[1.    0.72 0.72]
[[50   0   0]
 [ 0 36 14]
 [ 0 14 36]]
```

```
0.1
[1.          0.90196078 0.89795918]
[[50   0   0]
 [ 0 46   4]
 [ 0   6 44]]
```

```
1
[1.    0.92 0.92]
[[50   0   0]
 [ 0 46   4]
 [ 0   4 46]]
```

```
10
[0.94736842 0.88888889 0.8490566 ]
[[45   0   5]
 [ 0 44   6]
 [ 0   5 45]]
```

```
100
[0.41269841 0.04477612 0.11650485]
[[13 37   0]
 [ 0  3 47]
 [ 0 44   6]]
```


با افزایش مقدار Gamma از 0.01 به 1 عملکرد طبقه‌بند هم بهتر می‌شود اما با فاصله گرفتن زیاد آن از مقدار 1 (یعنی در Gamma= 100 و Gamma= 10) عملکرد طبقه‌بند کاهش می‌باید که علت آن (همانطور که در توضیحات پارامتر Gamma در بالا اشاره کردم) overfit شدن مدل ما به داده‌ی train می‌باشد.

C:

خروجی(ماتریس سطری از f1_scoreها که در سوال ۲ هم توضیح داده شد و ماتریس آشفته‌گی) به ازای مقادیر مختلف C. در خروجی ابتدا مقدار C چاپ شده و سپس f1_score کلاس‌ها و بعد هم ماتریس آشفته‌گی:

```
0.01
[0. 0. 0.]
[[ 0 45  5]
 [12  0 38]
 [ 5 45  0]]
```

```
0.1
[1.          0.48076923 0.4375    ]
[[50  0  0]
 [ 0 25 25]
 [ 0 29 21]]
```

```
1
[1.          0.90196078 0.89795918]
[[50  0  0]
 [ 0 46  4]
 [ 0  6 44]]
```

```
10
[1.          0.93069307 0.92929293]
[[50  0  0]
 [ 0 47  3]
 [ 0  4 46]]
```

```
100
[1.          0.92783505 0.93203883]
[[50  0  0]
 [ 0 45  5]
 [ 0  2 48]]
```

با افزایش مقدار C از 0.01 به 10 عملکرد طبقه‌بند هم بهتر می‌شود اما با فاصله گرفتن زیاد آن از مقدار 10 (یعنی در $C=100$) عملکرد طبقه‌بند کاهش می‌باید که علت آن (همانطور که در توضیحات پارامتر C در بالا اشاره کردم) **overfit** شدن مدل ما به داده‌ی **train** می‌باشد.

ج) نتیجه‌ی **grid search** با استفاده از 10-fold به صورت زیر می‌باشد:

```
Best parameters from gridsearch: {'C': 1, 'gamma': 0.1, 'kernel': 'poly'}
CV score=0.980
{'C': 1, 'gamma': 0.1, 'kernel': 'poly'}
[1.          0.93877551 0.94117647]
[[50  0  0]
 [ 0 46  4]
 [ 0  2 48]]
```

بهترین مقادیر گزارش شده توسط **grid search** برای پارامترها و همچنین نوع **kernel** در خروجی بالا قابل مشاهده است. می‌بینیم که با انتخاب **polynomial kernel** و مقدار $C = 1$ و $\text{Gamma} = 0.1$ برای پارامترها بیشترین عملکرد مدل محقق می‌شود. همانطور که در ماتریس آشفته‌گی مشخص است تنها ۶ داده به اشتباه طبقه‌بندی شده‌اند.

د) نتایج روش **one vs. all** برای تک تک کرنل‌ها به صورت زیر می‌باشد:

```
linear
[0.98989899 0.8125      0.83809524]
[[49  1  0]
 [ 0 39 11]
 [ 0  6 44]]
```

```
poly
[1.          0.94949495 0.95049505]
[[50  0  0]
 [ 0 47  3]
 [ 0  2 48]]
```

```
rbf
[1.          0.90384615 0.89583333]
[[50  0  0]
 [ 0 47  3]
 [ 0  7 43]]
```

می‌بینیم که با روش one vs. all، کرنل polynomial بهترین عملکرد را دارد.

نتایج روش one vs. one برای تک تک کرنل‌ها به صورت زیر می‌باشد:

```
linear
[1.    0.94 0.94]
[[50  0  0]
 [ 0 47  3]
 [ 0  3 47]]
.....
poly
[1.          0.92929293 0.93069307]
[[50  0  0]
 [ 0 46  4]
 [ 0  3 47]]
.....
rbf
[1.          0.90196078 0.89795918]
[[50  0  0]
 [ 0 46  4]
 [ 0  6 44]]
.....
```

می‌بینیم که با روش one vs. one، کرنل linear بهترین عملکرد را دارد.

همانطور که در پاسخ سوالات تشریحی اشاره کردم، امکان ایجاد عدم تقارن در داده‌های کلاس‌های اقلیت (کلاس‌های با داده‌ی کم)، توسط one vs. one classifier خیلی کم است. در نتیجه انتظار می‌رود که این روش عملکرد بهتری نسبت به روش one vs. all داشته باشد (هر چند کندتر است).