



به نام خدا

عنوان: پیاده‌سازی الگوریتم SOSStream

درس: داده کاوی

استاد درس: مهندس قاسمی

مهدی دهقانی

۲۲۰۷۹۷۰۴۸

mehdi.dehghani@ut.ac.ir

mahdiazadi18@yahoo.com

کلاس MicroCluster

این کلاس نشان دهنده خوشه‌های برنامه است که ویژگی‌هایی مثل تعداد داده‌ها، شعاع، مرکز، زمان ساخت، آخرین زمان تغییر و لیست داده‌های درون آن دارد.

این کلاس از هفت متد تشکیل شده است:

۱. `update_last_edited_time`

این متد آخرین زمان تغییر خوشه را به روز رسانی می‌کند.

۲. `insert`

این متد داده جدید را به خوشه اضافه می‌کند.

۳. `merge_data_points`

این متد برای ادغام دو لیست داده به کار می‌رود.

۴. `fading`

این متد برای محاسبه $f(t)$ که مقدار از بین رفتن خوشه است به کار می‌رود.

۵. `get_radius`

متد گرفتن مقدار شعاع خوشه

۶. `set_radius`

متد برای مقداردهی یا تغییر شعاع خوشه

۷. `get_centroid`

متد گرفتن مرکز خوشه

کلاس SOSTream

این کلاس الگوریتم SOSTream را پیاده‌سازی می‌کند و شامل متدهایی است که شبه کدهای آن‌ها را در فایل مقاله دیده‌ایم

این کلاس شامل متدهای زیر می‌باشد:

۱. find_neighbors

این متد برای پیدا کردن همسایه‌های خوشه برنده است.

۲. find_overlap

این متد برای پیدا کردن همسایه‌های خوشه خوشه برنده که با آن همپوشانی دارند به کار می‌رود.

۳. merge_clusters

این متد برای ادغام خوشه برنده با خوشه‌هایی که با آن همپوشانی دارند به کار می‌رود.

۴. update_cluster

به روز رسانی خوشه برنده و همسایگان آن هنگام اضافه شدن داده جدید

۵. fading_all

این متد برای از بین بردن خوشه‌ها با استفاده از محاسبه مقدار $f(t)$ و مقایسه آن با fade_threshold به کار می‌رود.

۶. adjust_centroid

این متد برای تنظیم مرکز خوشه به کار می‌رود که در متد `update_cluster` نیز استفاده شده است.

۷. insert

این متد برای درج یک خوشه در لیست خوشه‌ها به کار می‌رود.

۸. get_centroids_of_clusters

این متد برای گرفتن لیست مراکز خوشه‌ها به کار می‌رود.

۹. calculate_purity

این متد برای محاسبه خلوص خوشه‌ها و داده‌ها در این الگوریتم به کار می‌رود.

۱۰. باقی متدها

بقیه متدها برای گرفتن مقادیر مورد استفاده در الگوریتم هستند.

این اسکریپت برای اجرای الگوریتم SOSStream استفاده می‌شود. ابتدا دیتای مورد نظر را از فایل مربوطه می‌خوانیم و یک جایگشت رندوم روی آن اجرا می‌کنیم تا در زمان از بین بردن خوشه‌ها که نسبت به زمان ساخت و اضافه شدن خوشه‌ها و داده‌ها کار می‌کند به مشکلی بر نخوریم.

سپس مقادیر مورد نیاز، نظیر `minPts`، `merge_threshold`، `fade_threshold` و `لندا` را نسبت به شرایط مختلف مقداردهی می‌کنیم و الگوریتم را به صورت زیر اجرا می‌کنیم.

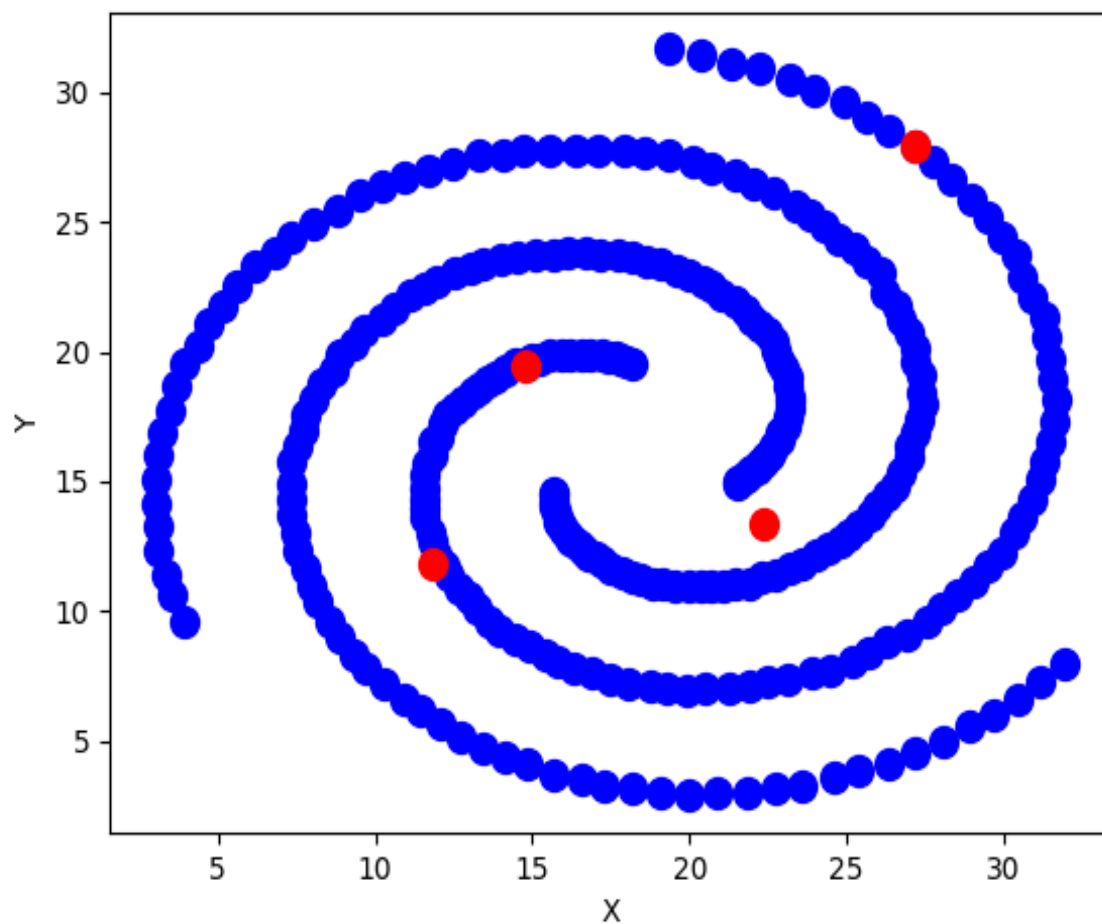
هر داده موجود در دیتاست را بصورت تکی به برنامه وارد می‌کنیم و نسبت به حالت فعلی داده‌ها و خوشه‌ها، نسبت به به روز رسانی یا ادغام یا موارد دیگر اقدام می‌کنیم. هر ۲۰ مرحله (هر بار که ۲۰ داده به برنامه وارد شده باشد) نیز از بین بردن خوشه‌ها را طبق `fading_all` انجام می‌دهیم و همچنین هر ۲۵ مرحله نیز تعداد خوشه‌ها و مقدار خلوص را ذخیره می‌کنیم تا در انتها در نتیجه کار تغییرات آن‌ها را ببینیم.

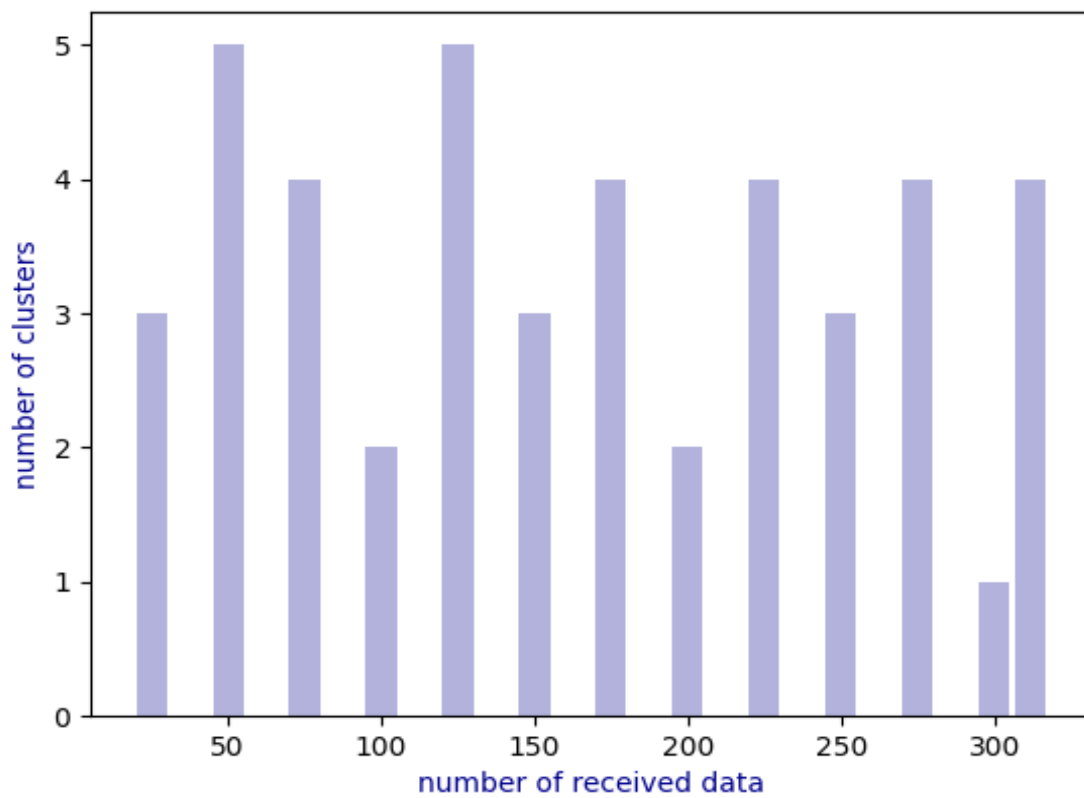
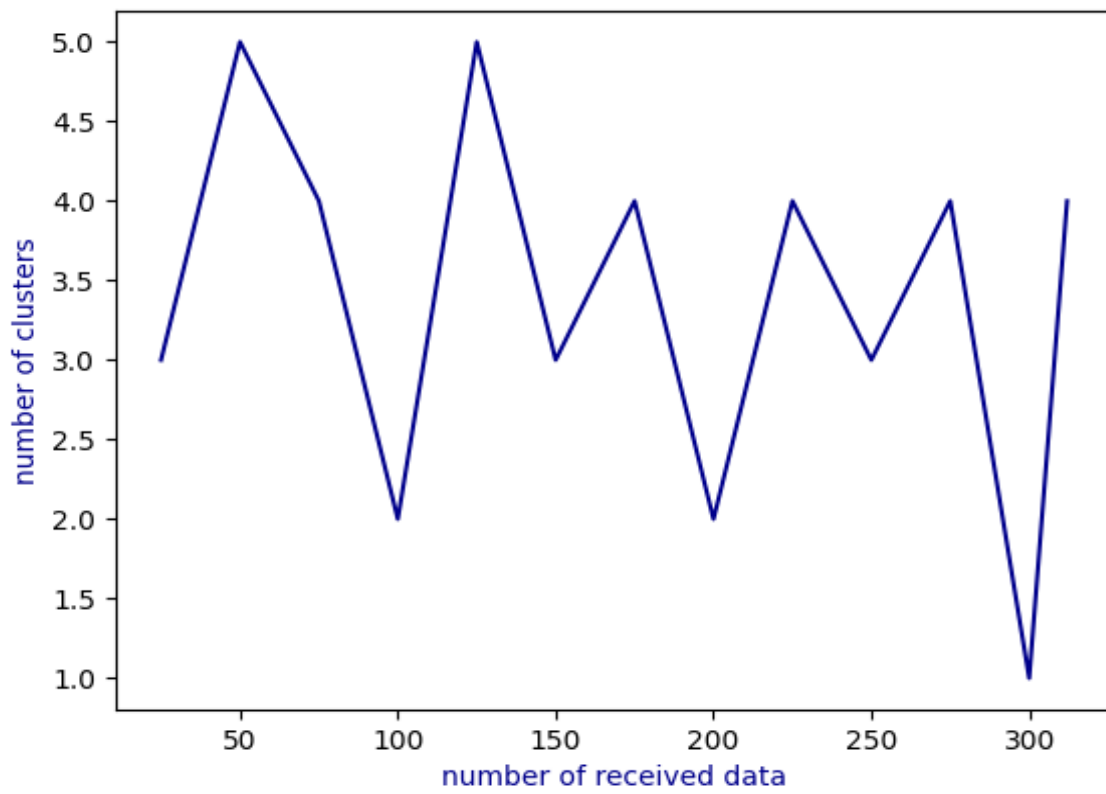
با توجه به محاسبه شدن موارد مختلف در الگوریتم و ذخیره کردن تعداد خوشه‌ها و مقدار خلوص در هر ۲۵ مرحله، می‌توانیم نمودارهای خطی و میله‌ای متناسب با آن‌ها را نمایش بدهیم و همچنین بعد از آن، نتایجی مثل زمان پردازش، تعداد خوشه‌های نهایی، تعداد خوشه‌های از بین رفته، تعداد خوشه‌های ادغام شده و میانگین خلوص را نشان داده و بررسی می‌کنیم.

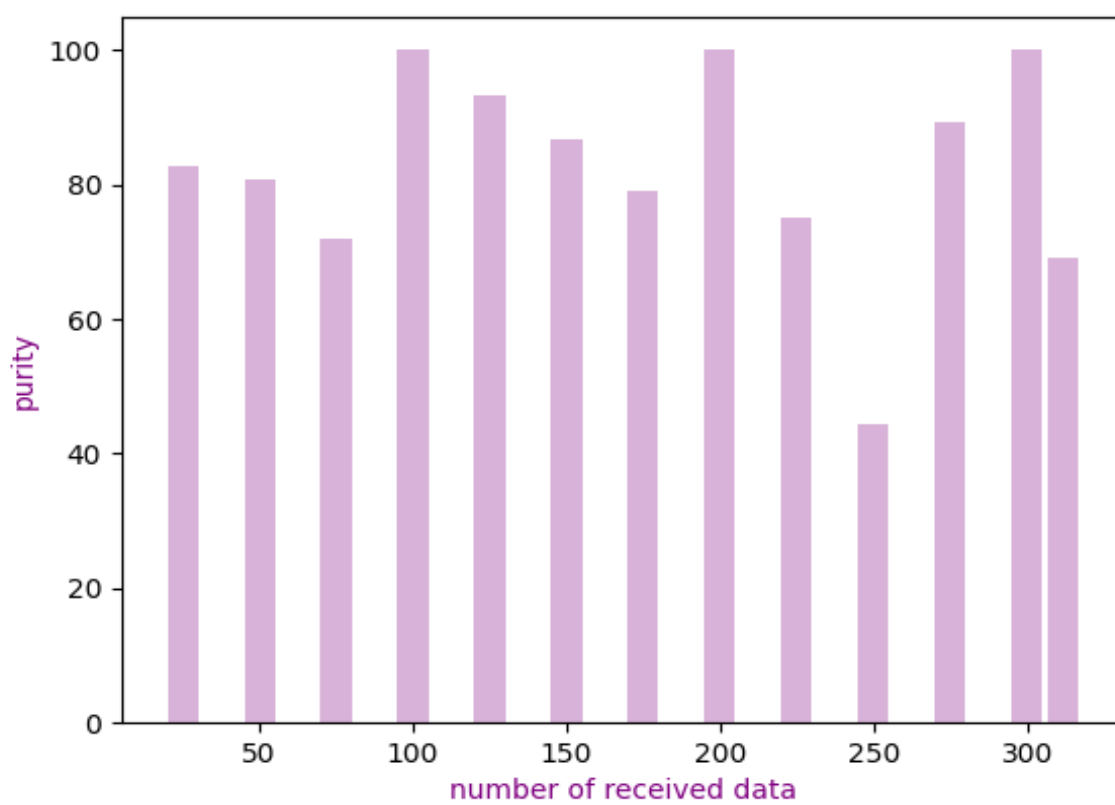
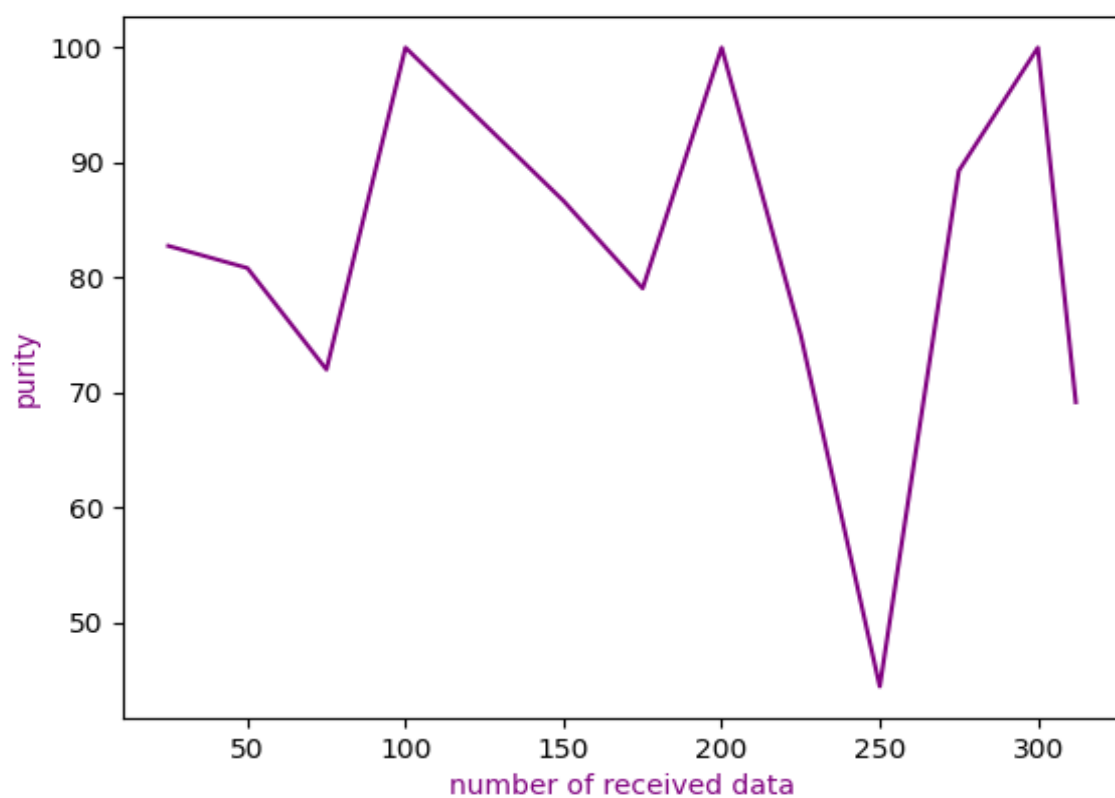
در ادامه مثال هایی از اجرا برنامه روی دو دیتاست قرار داده شده را می‌بینیم. بدیهی است که با تغییر مقادیر `آلفا`، `لندا`، `merge_threshold` یا ...، خروجی متفاوت خواهد شد. همچنین بخاطر جایگشت رندومی که در ابتدای کار روی دیتاست اعمال می‌شود نیز هر بار خروجی متفاوت خواهد بود.

در نمودارهای اول نقاط آبی رنگ، داده‌های ورودی و نقاط قرمز رنگ، مراکز خوشه‌های نهایی هستند.

```
# Create SOSTream object
sostream = SOSTream(data=data_frame,
                    alpha=0.1,
                    min_pts=3,
                    merge_threshold=7.75,
                    lambda_t=0.2,
                    fade_threshold=6.5
                    )
```



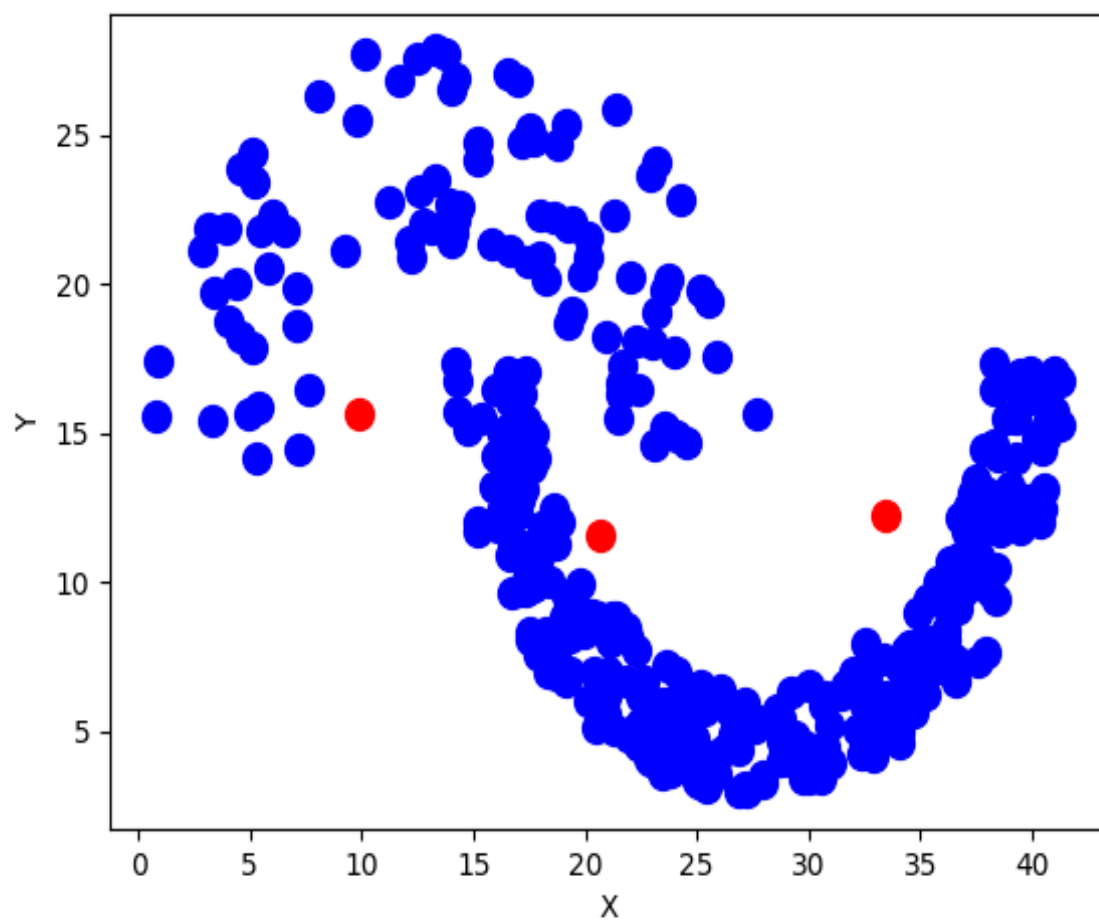


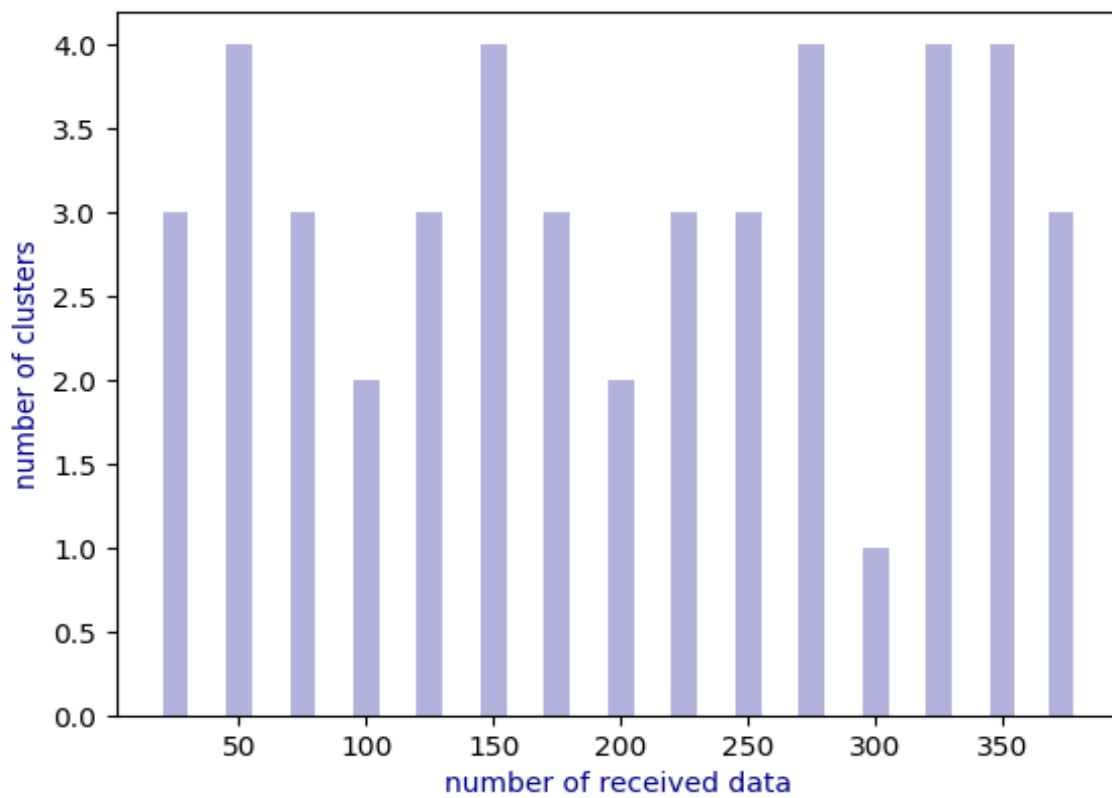
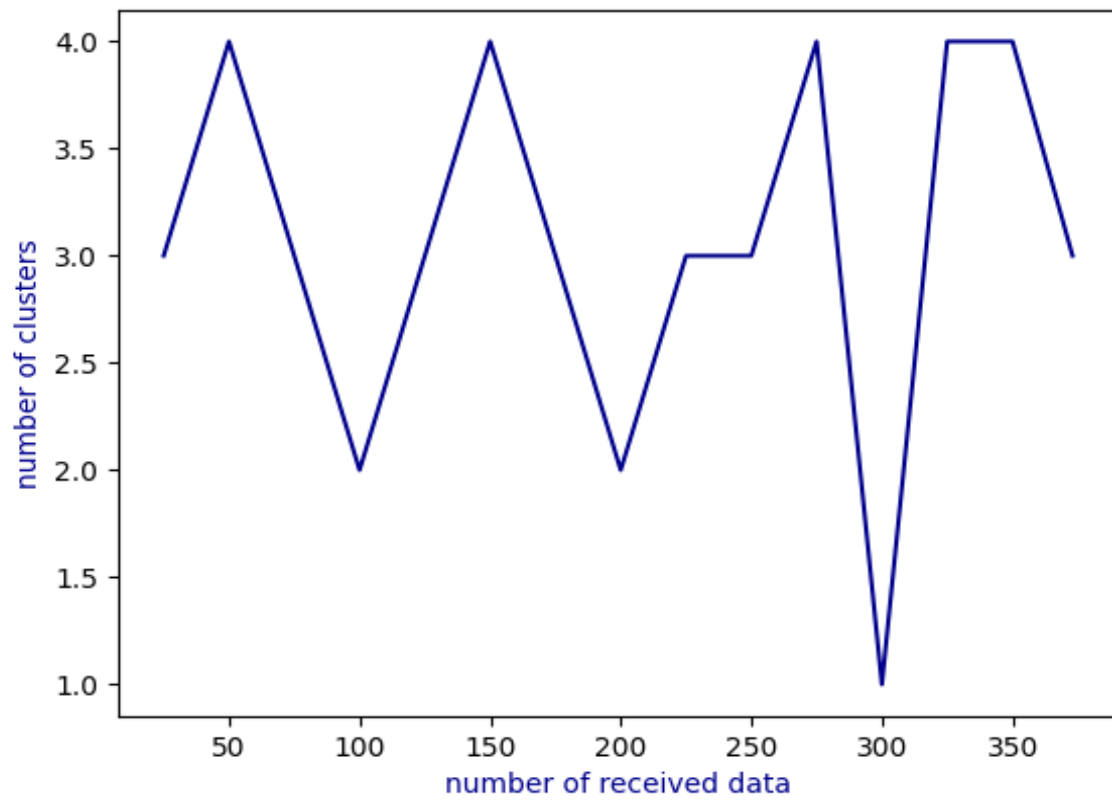


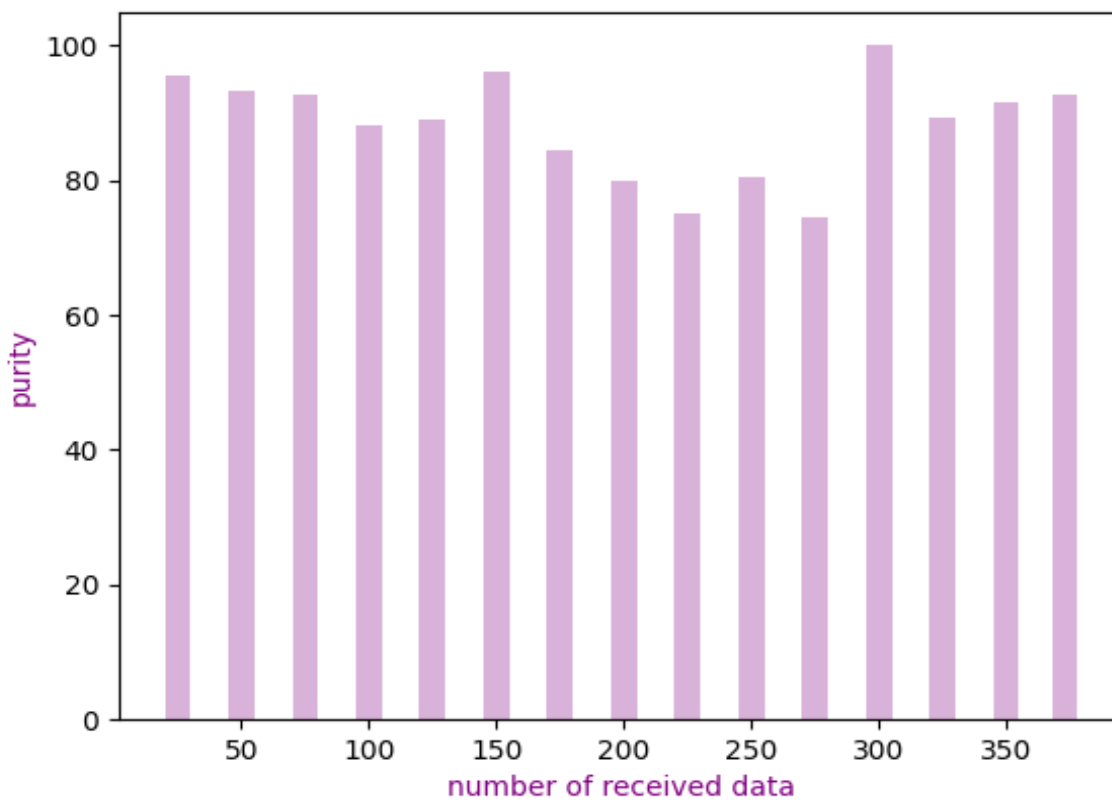
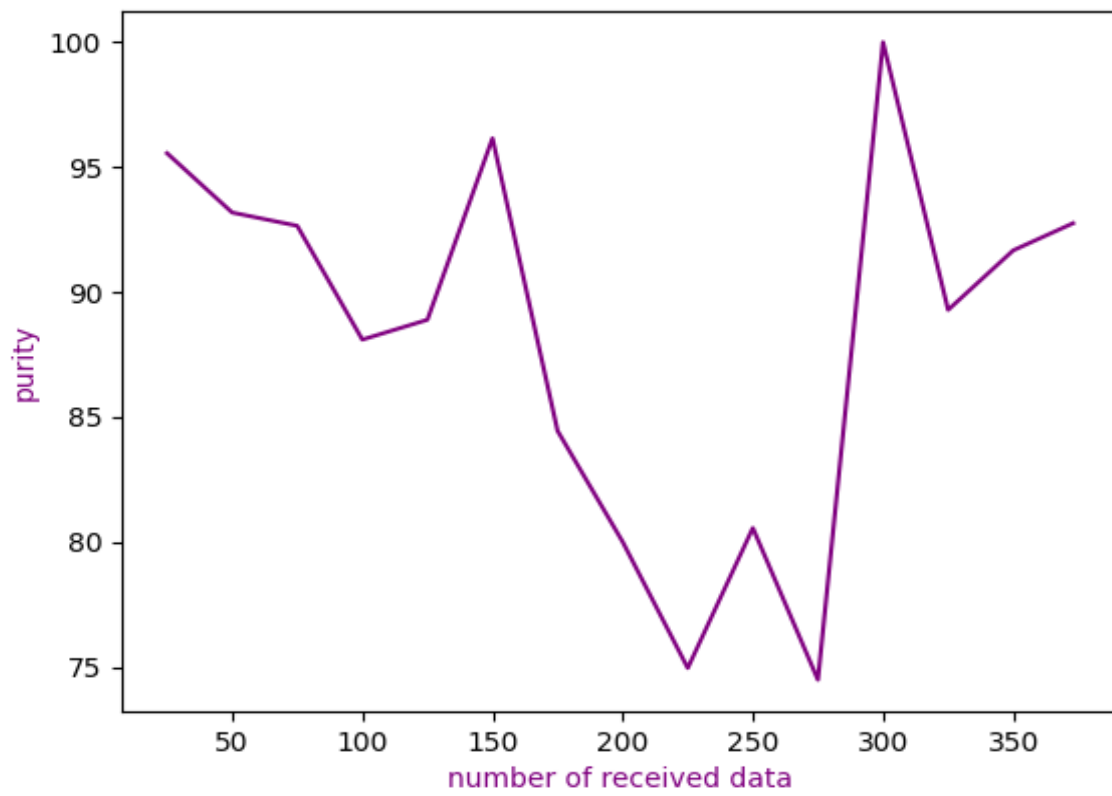
```
Process time =: 1.99  
Final number of clusters: 4  
Number of faded clusters: 31  
Number of merged clusters: 156  
Average purity: 82.50
```



```
# Create SOSTream object
sostream = SOSTream(data=data_frame,
                    alpha=0.1,
                    min_pts=3,
                    merge_threshold=9,
                    lambda_t=0.3,
                    fade_threshold=10
                    )
```







```
Process time =:  2.19  
Final number of clusters:  3  
Number of faded clusters:  37  
Number of merged clusters:  205  
Average purity:  88.18
```