# CSCD01

# Assignment 4

April 4th, 2021

## Table of Contents:

Bugs examined

**Issue #1: "Bisecting K-means #14214**
https://github.com/scikit-learn/scikit-learn/issues/14214

Issue Description:
Issue #14214 is a feature request for a clustering algorithm called Bisecting K-Means. Clustering algorithms are unsupervised learning algorithms that specialize in finding subgroups (called clusters) within datasets. In the scikit-learn library, there is already a similar algorithm implemented (called K-Means). However, Bisecting K-Means offers some things K-Means does not; K-Means typically identifies spherical shaped clusters only, while Bisecting K-Means does not have this limitation. Bisecting K-Means also offers better entropy measurement.
As a brief summary, K-Means divides a dataset into K clusters, each cluster having its own centroid (the mean of each cluster). K-Means starts by choosing K centroids, and improves upon estimation of the set of centroids.
On the other hand, Bisecting K-Means starts out by choosing two centroids of the whole dataset. It then chooses the cluster with the highest sum of squared errors (SSE), and replaces the centroid of such cluster with two more centroids by performing 2-Means (K-Means with 2 clusters). As such, conceptually, Bisecting K-Means uses K-Means within its own algorithm.
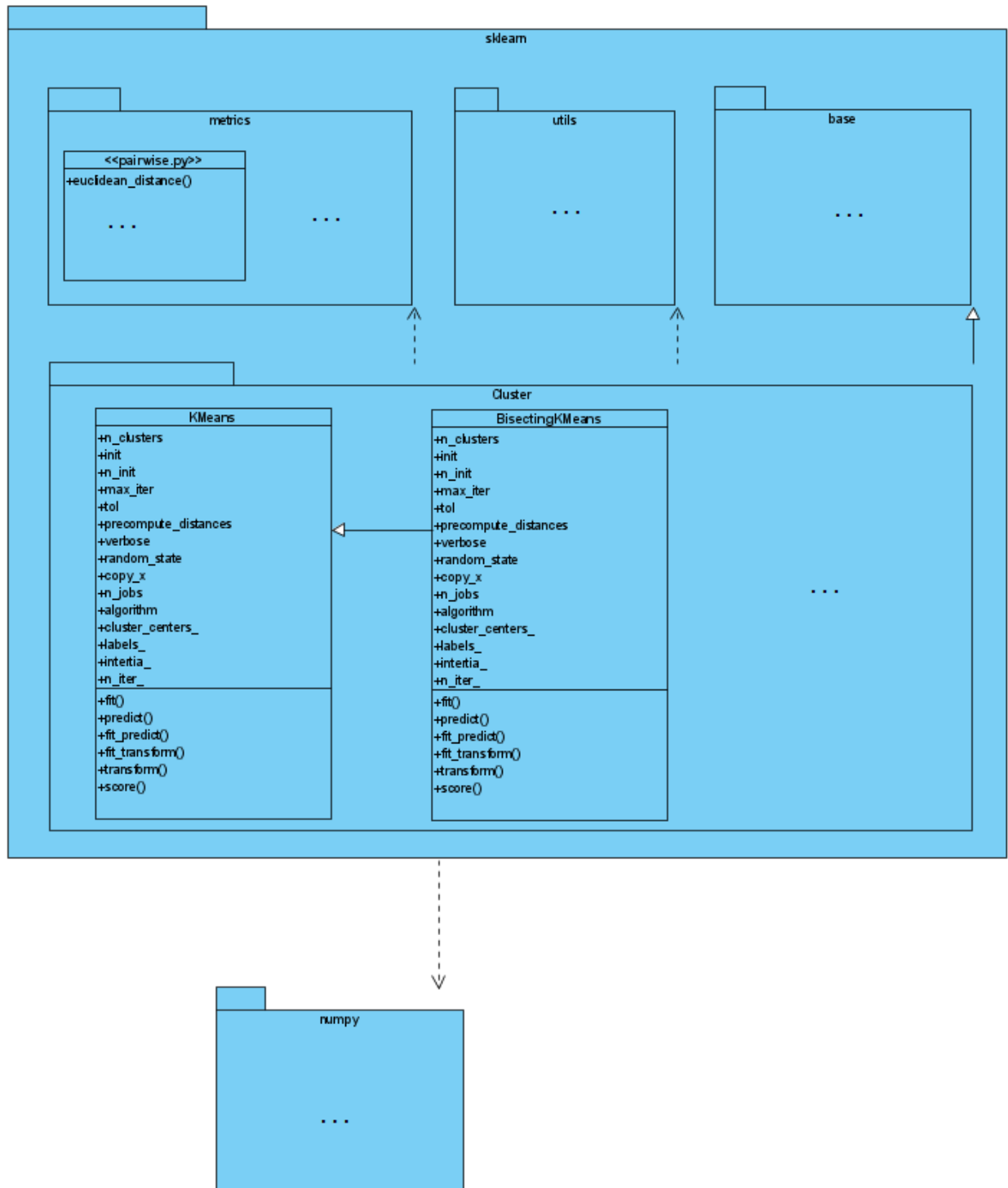
Issue investigation and design:
The changes to the library will be local to the clustering module, located at the path `sklearn/cluster/_kmeans.py`.
Each clustering algorithm in the library is its own class. Therefore, to stay consistent with the API, Bisecting K-Means will be its own class, much like Mini-Batch K-Means (another clustering algorithm related to K-Means) is its own class. Although the fit and predict methods are different, Mini-Batch K-Means inherits from K-Means. This makes sense because their instantiations are identical. This is the same case with Bisecting K-Means and K-Means, and as such, Bisecting K-Means will also inherit from K-Means. This also stays consistent with the architecture of the module as-is.
Since Bisecting K-Means uses K-Means within its own algorithm, it will also depend on K-Means in addition to inheritance. We would show this dependence by creating the K-Means object during the call to the `fit` method.
The change made to the module to incorporate this design is described by the following UML class diagram:

**sklearn**

**metrics**

<<pairwise.py>>
+euclidean_distance()

. . .

. . .

**utils**

. . .

**base**

. . .

**Cluster**

| KMeans |
| --- |
| +n_clusters |
| +init |
| +n_init |
| +max_iter |
| +tol |
| +precompute_distances |
| +verbose |
| +random_state |
| +copy_x |
| +n_jobs |
| +algorithm |
| +cluster_centers_ |
| +labels_ |
| +intertia_ |
| +n_iter_ |
| +fit() |
| +predict() |
| +fit_predict() |
| +fit_transform() |
| +transform() |
| +score() |

| BisectingKMeans |
| --- |
| +n_clusters |
| +init |
| +n_init |
| +max_iter |
| +tol |
| +precompute_distances |
| +verbose |
| +random_state |
| +copy_x |
| +n_jobs |
| +algorithm |
| +cluster_centers_ |
| +labels_ |
| +intertia_ |
| +n_iter_ |
| +fit() |
| +predict() |
| +fit_predict() |
| +fit_transform() |
| +transform() |
| +score() |

. . .

**numpy**

. . .

All of the attributes will be identical, as most will be passed on to KMeans when the algorithm needs to perform 2-Means clustering. An exception would be `n_clusters`, as this will dictate the final number of clusters to end on; Bisecting KMeans will pass on the value of "2" to KMeans in every step that requires 2-Means clustering.

In terms of the algorithm; the implementation will heavily rely on numpy arrays and the operations offered by numpy arrays, as well as the `eculidean_distance()` function offered by the metrics module of sklearn.
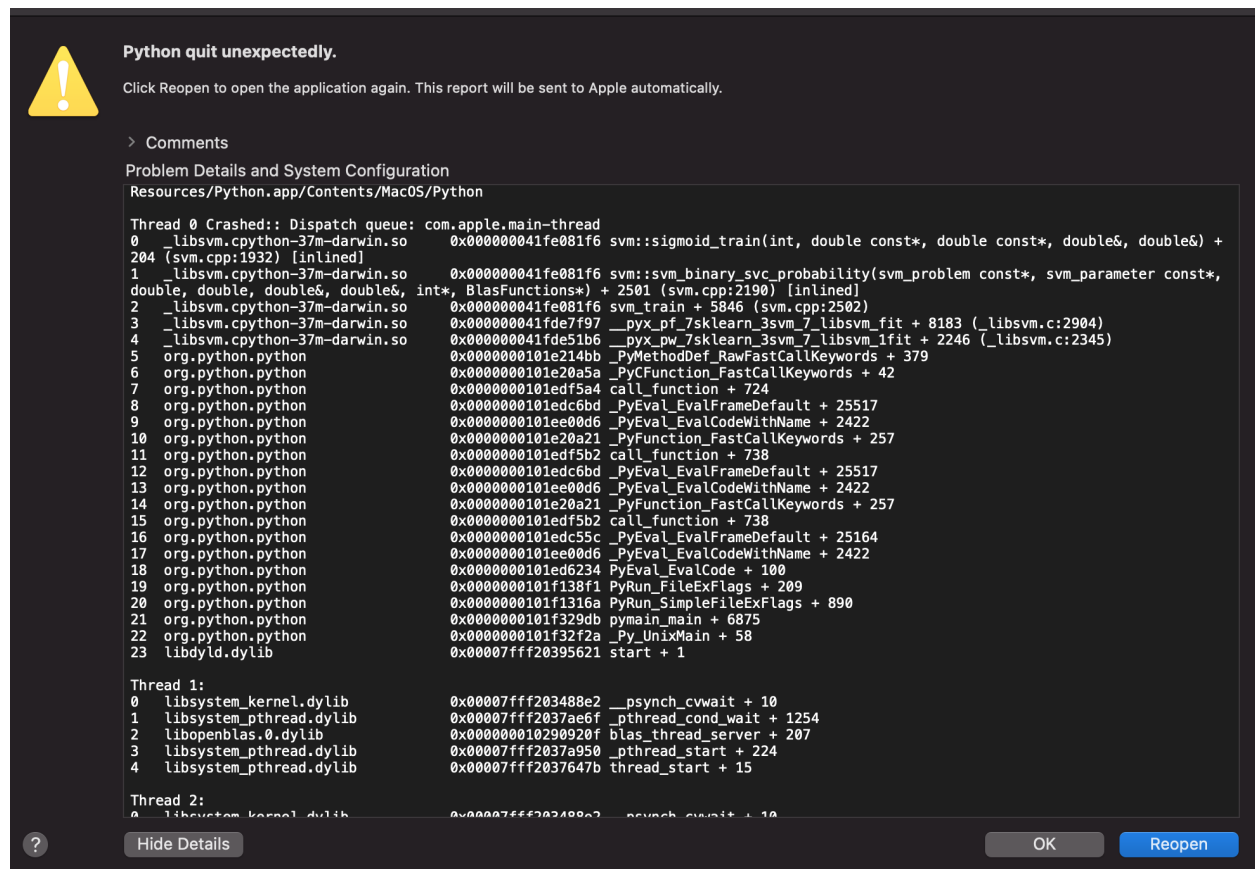
**Issue #2: "Segmentation Fault with SVC with Probability = True and degree 3"**
https://github.com/scikit-learn/scikit-learn/issues/15008

Issue Description:
A segmentation fault occurs when attempting to instantiate a SVC object with probability = True and degree 3. When running the code below, users will hit a segmentation fault after about ten minutes.

```
import numpy as np
features = np.random.rand(2100000,768)
target = np.random.randint(low =100000,high=500000, size=2100000)
from sklearn.svm import SVC
svc =
SVC(kernel='poly',gamma=100,C=10,degree=3,decision_function_shape='ovo',probability=Tr
ue,random_state=42,verbose=True)
svc.fit(features,target)
```

**Implementation and Fix**:

The problem occurs in the file *sklearn/svm/src/libsvm/svm.cpp*, specifically in the function *sigmoid_train*. On first glance it looks like the issue is with an index overflowing as suggested by a commenter on the issue. However, after debugging, the error is with a malloc in the function *\*PREFIX(train)*. When the variable nr_classes gets above ~390000 the error occurs since it mallocs something of size 79154844021 (nr-classes\*(nr_classes-1)/2). This is above the maximum value for an integer in C++, (see figure 2) and thus a segmentation fault occurs.

By casting this to an int64_t the segmentation vault is avoided, as by figure 1, the maximum value for an int64_t is 9,223,372,036,854,775,807, which is greater than 79154844021.
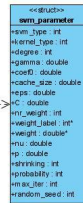
As a preventative measure we also looked at other areas of the code where this integer overflow could happen. To ensure the error doesn't also occur there, we changed the appropriate variables to an int64_t.

|  | int | int64_t |
|---|---|---|
| Memory allocation size | 4 bytes | 8 bytes |
| Max Value | 2,3147,483,647 | 9,223,372,036,854,775,807 |

*Figure 1*

**UML Changes**

5

The bug in question concerns/affects the svm section of the UML as seen here:

svm_csr

**ONE_CLASS_Q**
+__pad0__
+x : prob
+param : prob
+bias_functions : prob
-QD : double
-cache : Cache*
+for()
+get_Q()
+get_QD()
+swap_index()
+~ONE_CLASS_Q()

<<struct>>
**decision_function**
+alpha : double*
+rho : double

**Kernel**
#kernel_function(int i, int j) const *
-x_square : double*
-m_blas : BlasFunctions*
-kernel_type : const int
-degree : const int
-gamma : const double
-coef0 : const double
+Kernel()
+get_Q()
+get_QD()
+swap_index()
+is_function()
-PREFIX()
-kernel_linear()
-kernel_poly()
-kernel_rbf()
-kernel_sigmoid()
-kernel_precomputed()
-dot()

**SVC_Q**
+__pad0__
+x : prob
+param : prob
+bias_functions : prob
-cache : Cache*
-QD : double
-y : schar*
+for()
+get_Q()
+get_QD()
+swap_index()
+~SVC_Q()

-cache
**Cache**
-l : int
-size : long int
-head : head_t*
-lru_head : head_t
+Cache()
+~Cache()
+get_data()
+swap_index()
-lru_delete()
-lru_insert()
-lru_delete()
-lru_insert()

<<struct>>
**head_t**
+prev : head_t*
+next : head_t*
-data : Qfloat*
+len : int

<<struct>>
<<Typedef>>
**BlasFunctions**
+dot : dot_func

-bias

<<enumeration>>
**UNNAMED_ENUM**
LOWER_BOUND
UPPER_BOUND
FREE

**SVR_Q**
+__pad0__
+x : prob
+param : prob
+bias_functions : prob
-cache : Cache*
-QD : double
-sign : schar
-index : int
-buffer : Qfloat[0]
-next_buffer : int
-l : int
+for()
+swap_index()
+get_Q()
+get_QD()
+~SVR_Q()

<<Typedef>>
**Qfloat**
-data

**Solver**
#active_size : int
-y : schar*
#G : double*
#alpha_status : char*
#alpha : double*
-Q : QMatrix const*
#QD : double const*
#eps : double
#Cp : double
#Cn : double
#C : double*
#p : double*
#active_set : int*
#G_bar : double*
#l : int
#unshrink : bool
+Solver()
+~Solver()
+Solve()
#get_C()
#update_alpha_status()
#is_upper_bound()
#is_lower_bound()
#is_free()
#swap_index()
#reconstruct_gradient()
#select_working_set()
#calculate_rho()
#do_shrinking()
-be_shrunk()
+Solve()

<<struct>>
**SolutionInfo**
+obj : double
+rho : double
+upper_bound : double*
+r : double
+solve_timed_out : bool

**QMatrix**
+get_Q()
+get_QD()
+swap_index()
+~QMatrix()

**Solver_NU**
-si : SolutionInfo*
+Solver_NU()
+Solve()
-select_working_set()
-calculate_rho()
-be_shrunk()
-do_shrinking()
+Solve()

**GLOBAL**
+min()
+max()
+swap()
+clone()
+pow()
+print_string_stdout()
+info()
+remove_zero_weight()
+PREFIX()
+if()
+if()
+PREFIX()
+for()
+for()
+info()
+free()
+free()
+free()
+free()
+free()
+free()
+free()
+free()
+free()
+free()
+free()
+free()
+cross_validation()
+get_svm_type()
+get_nr_class()
+get_labels()
+get_svr_probability()
+predict_values()
+predict()
+predict_probability()
+free_model_content()
+free_and_destroy_model()
+destroy_param()
+check_parameter()
+set_print_string_function()
+svm_train()
+svm_cross_validation()
+svm_save_model()
+svm_load_model()
+svm_get_svm_type()
+svm_get_nr_class()
+svm_get_labels()
+svm_get_svr_probability()
+svm_predict_values()
+svm_predict()
+svm_predict_probability()
+svm_free_model_content()
+svm_free_and_destroy_model()
+svm_destroy_param()
+svm_check_parameter()
+svm_set_print_string_function()
+svm_csr_train()
+svm_csr_cross_validation()
+svm_csr_get_svm_type()
+svm_csr_get_nr_class()
+svm_csr_get_labels()
+svm_csr_get_svr_probability()
+svm_csr_predict_values()
+svm_csr_predict()
+svm_csr_predict_probability()
+svm_csr_free_model_content()
+svm_csr_free_and_destroy_model()
+svm_csr_destroy_param()
+svm_csr_check_parameter()

-y
<<Typedef>>
**schar**

<<struct>>
**svm_parameter**
+svm_type : int
+kernel_type : int
+degree : int
+gamma : double
+coef0 : double
+cache_size : double
+eps : double
+C : double
+nr_weight : int
+weight_label : int*
+weight : double*
+nu : double
+p : double
+shrinking : int
+probability : int
+max_iter : int
+random_seed : int

In terms of the UML impact for our bug, no new classes or functions were added to the uml. For the bug in question, the main part of the uml we changed and are concerned with is the svm_train() function as displayed here:

```
                GLOBAL
+min()
+max()
+swap()
+clone()
+powi()
+print_string_stdout()
+info()
+remove_zero_weight()
+PREFIX()
+if()
+if()
+PREFIX()
+for()
+for()
+info()
+free()
+free()
+free()
+free()
+free()
+free()
+free()
+free()
+free()
+free()
+free()
+free()
+free()
+free()
+cross_validation()
+get_svm_type()
+get_nr_class()
+get_labels()
+get_svr_probability()
+predict_values()
+predict()
+predict_probability()
+free_model_content()
+free_and_destroy_model()
+destroy_param()
+check_parameter()
+set_print_string_function()
+svm_train()                          <---
+svm_cross_validation()
+svm_save_model()
+svm_load_model()
+svm_get_svm_type()
+svm_get_nr_class()
+svm_get_labels()
+svm_get_svr_probability()
+svm_predict_values()
+svm_predict()
+svm_predict_probability()
+svm_free_model_content()
+svm_free_and_destroy_model()
+svm_destroy_param()
+svm_check_parameter()
+svm_set_print_string_function()
+svm_csr_train()
+svm_csr_cross_validation()
+svm_csr_get_svm_type()
+svm_csr_get_nr_class()
+svm_csr_get_labels()
+svm_csr_get_svr_probability()
+svm_csr_predict_values()
+svm_csr_predict()
+svm_csr_predict_probability()
+svm_csr_free_model_content()
+svm_csr_free_and_destroy_model()
+svm_csr_destroy_param()
+svm_csr_check_parameter()
```

**Bug Selected:**

For our bug fix, we selected bug #2 to implement. The choice was made because we felt it was an opportunity to make a significant contribution to a frequently used feature. We are also confident about the fix's chance to merge into the main sci-kit learn branch.

**Why This Bug is "Hard":**

Our chosen issue did not have a labelled difficulty. We determined that the fix should be considered "hard" for the following reasons:

1. The stack trace for the error is not actually where the bug is located
2. The code takes upwards of 10 minutes to run and produce the error
3. The fix was in C++, a language our group is less comfortable working in

We verified these reasons through email correspondences with our TA's, and were given approval to fix the bug.

**Testing:**

*Unit Testing:*

We do not believe this fix requires unit testing since it takes many hours to fully run, and uses a lot of computing resources. However, for the sake of the assignment, to unit test, we would create a test that runs the code above. If the process terminates without a segmentation fault, it would mean that our fix is successful.

*Acceptance Testing:*

A customer acceptance test to validate the correctness of our fix would go as follows:

1. Create a matrix represented as a nested array of size (2100000 x 768) and populate it with random samples of a uniform distribution over 0,1
2. Create an array of size 2100000, with random values in it between 100,000 and 500,000.
3. Instantiate an SVC object with specific parameters specified above
4. Try to fit the matrix and array to the SVC object, using the SVC's fit method

Our fix also passed scikit-learn's suite of acceptance tests, as shown here in our pending pull request: https://github.com/scikit-learn/scikit-learn/pull/19782