


```
In [1]: import scipy as sp
import scipy as sp
import scipy.cluster.vq
import scipy.spatial.distance
import scipy.stats
import sklearn.cluster

import pylab as pl

dst = sp.spatial.distance.euclidean

def gap_statistics(data, refs=None, nrefs=20, ks=range(1,11)):
    """Computes the gap statistics for an nxm dataset.

    The gap statistic measures the difference between within-cluster dispersion on an input
    dataset and that expected under an appropriate reference null distribution.

    Computation of the gap statistic, then, requires a series of reference (null) distributions.
    One may either input a precomputed set of reference distributions (via the parameter refs)
    or specify the number of reference distributions (via the parameter nrefs) for automatic
    generation of uniform distributions within the bounding box of the dataset (data).

    Each computation of the gap statistic requires the clustering of the input dataset and of
    several reference distributions. To identify the optimal number of clusters  $k$ , the gap
    statistic is computed over a range of possible values of  $k$  (via the parameter ks).

    For each value of  $k$ , within-cluster dispersion is calculated for the input dataset and each
    reference distribution. The calculation of the within-cluster dispersion for the reference
    distributions will have a degree of variation, which we measure by standard deviation or
    standard error.

    The estimated optimal number of clusters, then, is defined as the smallest value  $k$  such that
     $gap_k$  is greater than or equal to the sum of  $gap_{k+1}$  minus the expected error  $err_{k+1}$ .

    Args:
        data ((n,m) SciPy array): The dataset on which to compute the gap statistics.
        refs ((n,m,k) SciPy array, optional): A precomputed set of reference distributions.
            Defaults to None.
        nrefs (int, optional): The number of reference distributions for automatic generation.
            Defaults to 20.
        ks (list, optional): The list of values  $k$  for which to compute the gap statistics.
            Defaults to range(1,11), which creates a list of values from 1 to 10.
```

Returns:

gaps: an array of gap statistics computed for each k.

errs: an array of standard errors (se), with one corresponding to each gap computation.

difs: an array of differences between each gap_k and the sum of gap_{k+1} minus err_{k+1}.

"""

shape = data.shape

if refs==None:

 tops = data.max(axis=0) # maxima along the first axis (rows)

 bots = data.min(axis=0) # minima along the first axis (rows)

 dists = sp.matrix(sp.diag(tops-bots)) # the bounding box of the input dataset

 # Generate nrefs uniform distributions each in the half-open interval [0.0, 1.0)

 rands = sp.random.random_sample(size=(shape[0],shape[1], nrefs))

 # Adjust each of the uniform distributions to the bounding box of the input dataset

 for i in range(nrefs):

 rands[:, :, i] = rands[:, :, i]*dists+bots

else:

 rands = refs

gaps = sp.zeros((len(ks),)) # array for gap statistics (length ks)

errs = sp.zeros((len(ks),)) # array for model standard errors (length ks)

difs = sp.zeros((len(ks)-1,)) # array for differences between gaps (length ks-1)

for (i,k) in enumerate(ks): # iterate over the range of k values

 # Cluster the input dataset via k-means clustering using the current value of k

 try:

 (kmc,kml) = sp.cluster.vq.kmeans2(data, k)

 except LinAlgError:

 kmeans = sklearn.cluster.KMeans(n_clusters=k).fit(data)

 (kmc, kml) = kmeans.cluster_centers_, kmeans.labels_

 # Generate within-dispersion measure for the clustering of the input dataset

 disp = sum([dst(data[m,:],kmc[kml[m],:]) for m in range(shape[0])])

 # Generate within-dispersion measures for the clusterings of the reference datasets

 refdisps = sp.zeros((rands.shape[2],))

 for j in range(rands.shape[2]):

 # Cluster the reference dataset via k-means clustering using the current value of k

 try:

```

        (kmc, kml) = sp.cluster.vq.kmeans2(rands[:, :, j], k)
    except LinAlgError:
        kmeans = sklearn.cluster.KMeans(n_clusters=k).fit(rands[:, :, j])
        (kmc, kml) = kmeans.cluster_centers_, kmeans.labels_

    refdisps[j] = sum([dst(rands[m, :, j], kmc[kml[m], :]) for m in range(shape[0])])

    # Compute the (estimated) gap statistic for k
    gaps[i] = sp.mean(sp.log(refdisps) - sp.log(disps))

    # Compute the expected error for k
    errs[i] = sp.sqrt(sum(((sp.log(refdisp) - sp.mean(sp.log(refdisps)))**2) \
                          for refdisp in refdisps)/float(nrefs)) * sp.sqrt(1+1/nrefs))

    # Compute the difference between gap_k and the sum of gap_k+1 minus err_k+1
    difs = sp.array([gaps[k] - (gaps[k+1] - errs[k+1]) for k in range(len(gaps)-1)])

    #print "Gaps: " + str(gaps)
    #print "Errs: " + str(errs)
    #print "Difs: " + str(difs)

    return gaps, errs, difs

def plot_gap_statistics(gaps, errs, difs):
    """Generates and shows plots for the gap statistics.

    A figure with two subplots is generated. The first subplot is an errorbar plot of the
    estimated gap statistics computed for each value of k. The second subplot is a barplot
    of the differences in the computed gap statistics.

    Args:
        gaps (SciPy array): An array of gap statistics, one computed for each k.
        errs (SciPy array): An array of standard errors (se), with one corresponding to each gap
            computation.
        difs (SciPy array): An array of differences between each gap_k and the sum of gap_k+1
            minus err_k+1.

    """
    # Create a figure
    fig = pl.figure(figsize=(16, 4))

    pl.subplots_adjust(wspace=0.35) # adjust the distance between figures

```

```
# Subplot 1
ax = fig.add_subplot(121)
ind = range(1, len(gaps)+1) # the x values for the gaps

# Create an errorbar plot
rects = ax.errorbar(ind, gaps, yerr=errs, xerr=None, linewidth=1.0)

# Add figure labels and ticks
ax.set_title('Clustering Gap Statistics', fontsize=16)
ax.set_xlabel('Number of clusters k', fontsize=14)
ax.set_ylabel('Gap Statistic', fontsize=14)
ax.set_xticks(ind)

# Add figure bounds
ax.set_ylim(0, max(gaps+errs)*1.1)
ax.set_xlim(0, len(gaps)+1.0)

# Subplot 2
ax = fig.add_subplot(122)
ind = range(1, len(difs)+1) # the x values for the difs

max_gap = None
if len(np.where(difs > 0)[0]) > 0:
    max_gap = np.where(difs > 0)[0][0] + 1 # the k with the first positive dif

# Create a bar plot
ax.bar(ind, difs, alpha=0.5, color='g', align='center')

# Add figure labels and ticks
if max_gap:
    ax.set_title('Clustering Gap Differences\n(k=%d Estimated as Optimal)' % (max_gap), \
                fontsize=16)
else:
    ax.set_title('Clustering Gap Differences\n', fontsize=16)
ax.set_xlabel('Number of clusters k', fontsize=14)
ax.set_ylabel('Gap Difference', fontsize=14)
ax.xaxis.set_ticks(range(1, len(difs)+1))

# Add figure bounds
ax.set_ylim(min(difs)*1.2, max(difs)*1.2)
ax.set_xlim(0, len(difs)+1.0)

# Show the figure
```

```
pl.show()
```

```

In [3]: # (c) 2014 Reid Johnson
# BSD License
#
# Function to compute the sum of squared distance (SSQ) for evaluating k-means clustering.

import numpy as np
import scipy as sp
import sklearn.cluster
from scipy.spatial.distance import cdist, pdist

import pylab as pl

def ssq_statistics(data, ks=range(1,11), ssq_norm=True):
    """Computes the sum of squares for an nxm dataset.

    The sum of squares (SSQ) is a measure of within-cluster variation that measures the sum of
    squared distances from cluster prototypes.

    Each computation of the SSQ requires the clustering of the input dataset. To identify the
    optimal number of clusters k, the SSQ is computed over a range of possible values of k
    (via the parameter ks). For each value of k, within-cluster dispersion is calculated for the
    input dataset.

    The estimated optimal number of clusters, then, is defined as the value of k prior to an
    "elbow" point in the plot of SSQ values.

    Args:
        data ((n,m) SciPy array): The dataset on which to compute the gap statistics.
        ks (list, optional): The list of values k for which to compute the gap statistics.
            Defaults to range(1,11), which creates a list of values from 1 to 10.

    Returns:
        ssqs: an array of SSQs, one computed for each k.

    """
    ssqs = sp.zeros((len(ks),)) # array for SSQs (length ks)

    #n_samples, n_features = data.shape # the number of rows (samples) and columns (features)
    #if n_samples >= 2500:
    #    # Generate a small sub-sample of the data
    #    data_sample = shuffle(data, random_state=0)[:1000]
    #else:

```

```

# data_sample = data

for (i,k) in enumerate(ks): # iterate over the range of k values
    # Fit the model on the data
    kmeans = sklearn.cluster.KMeans(n_clusters=k, random_state=0).fit(data)

    # Predict on the data (k-means) and get labels
    #labels = kmeans.predict(data)

    if ssq_norm:
        dist = np.min(cdist(data, kmeans.cluster_centers_, 'euclidean'), axis=1)

        tot_withinss = sum(dist**2) # Total within-cluster sum of squares
        totss = sum(pdist(data)**2) / data.shape[0] # The total sum of squares
        betweenss = totss - tot_withinss # The between-cluster sum of squares
        ssqs[i] = betweenss/totss*100
    else:
        # The sum of squared error (SSQ) for k
        ssqs[i] = kmeans.inertia_

return ssqs

def plot_ssq_statistics(ssqs):
    """Generates and shows plots for the sum of squares (SSQ).

    A figure with one plot is generated. The plot is a bar plot of the SSQ computed for each
    value of k.

    Args:
        ssqs (SciPy array): An array of SSQs, one computed for each k.

    """
    # Create a figure
    fig = pl.figure(figsize=(6.75, 4))

    ind = range(1,len(ssqs)+1) # the x values for the ssqs
    width = 0.5 # the width of the bars

    # Create a bar plot
    #rects = pl.bar(ind, ssqs, width)
    pl.plot(ind, ssqs)

    # Add figure labels and ticks

```



```
pl.title('Clustering Sum of Squared Distances', fontsize=16)
pl.xlabel('Number of clusters k', fontsize=14)
pl.ylabel('Sum of Squared Distance (SSQ)', fontsize=14)
pl.xticks(ind)

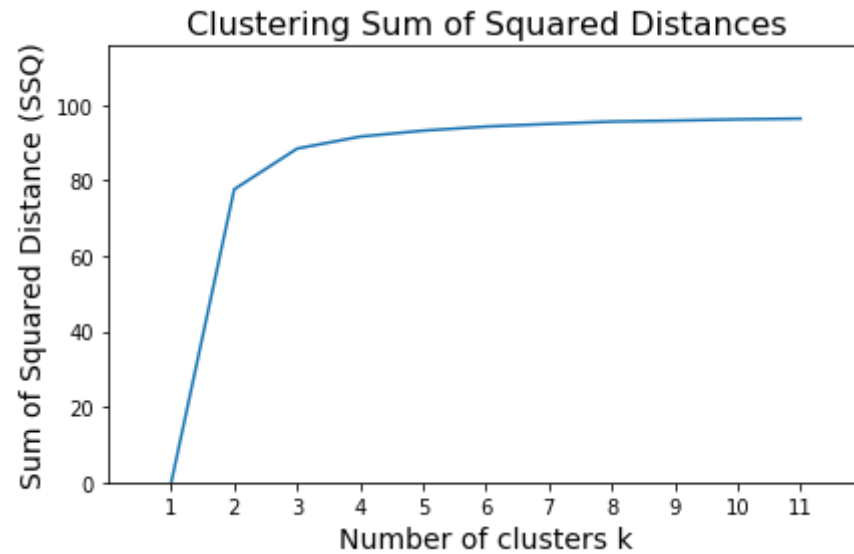
# Add text labels
#for rect in rects:
#    height = rect.get_height()
#    pl.text(rect.get_x()+rect.get_width()/2., 1.05*height, '%d' % int(height), \
#            ha='center', va='bottom')

# Add figure bounds
pl.ylim(0, max(ssqs)*1.2)
pl.xlim(0, len(ssqs)+1.0)

pl.show()
```

```
In [4]: from sklearn import datasets
# Load the Iris flower dataset
iris = datasets.load_iris()
data = iris.data
```

```
In [5]: # Generate and plot the SSQ statistics
ssqs = ssq_statistics( data , ks=range( 1 , 11+ 1 ) )
plot_ssq_statistics(ssqs)
```



```

In [61]: # Keep track of vals
maxGapList = [0]*10
# Python program to get average of a list
def average(lst):
    return sum(lst) / len(lst)
# get k values of 1-10 inclusive
for i in range(1,11):

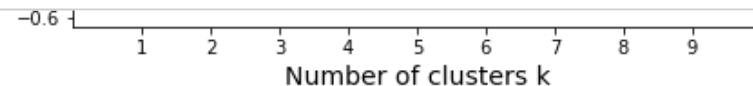
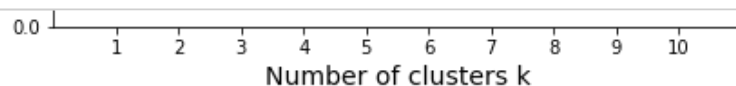
    print("SSQ and Gap Statistics for k value of: ",i)
    # Generate and plot the SSQ statistics
    ssqs = ssq_statistics(data, ks=range(1,10+1))
    plot_ssq_statistics(ssqs)

    # Generate and plot the gap statistics
    gaps, errs, difs = gap_statistics(data, nrefs=20, ks=range(1,10+1))
    plot_gap_statistics(gaps, errs, difs)

    # Initialize max gap to 0 each iteration
    maxGap = 0
    #Find first positive value and make predictions or else skip and reset max gap to 0
    try:
        len(np.where(difs > 0)[0]) > 0
        maxGap = np.where(difs > 0)[0][0] + 1
        print("Max Gap: ", maxGap, "\n")
        # Add to list tally
        maxGapList[i-1] = maxGap
        kmeans = sklearn.cluster.KMeans(n_clusters=maxGap, random_state=0).fit(data)
        # Make prediction of newest values
        labels = kmeans.predict(data)
    except:
        print("No valid max gap")

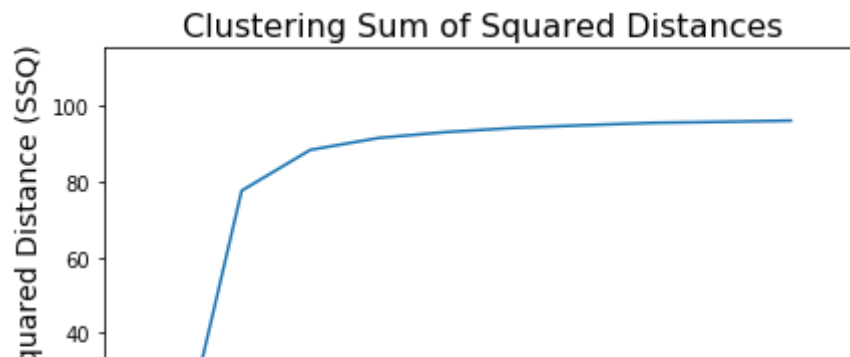
print("List of all current maxGaps", maxGapList, "\n")
avgK = average(maxGapList)
print("Average k: ", avgK, "\n")

```



Max Gap: 4

SSQ and Gap Statistics for k value of: 8



1. Where did you estimate the elbow point to be (between what values of k)? What value of k was typically estimated as optimal by the gap statistic? To adequately answer this question, consider generating both measures several times, as there may be some amount of variation in the value of k that they each estimate as optimal.

The elbow point can be estimated being between $k=2$ and $k=4$. After running the program a few times and observing the "Average K" output, the optimal k value is typically between 3 and 3.3 with a maximum of 5 and a minimum of 2.

2. How close are the estimates generated by the elbow point and gap statistic to the number of species of Iris represented in the dataset?

There are 3 known species of Iris in the dataset and with our estimates converging just around 3, this is a very accurate estimation. However we can see that there is some varying error of SSQ and gap statistics over iterations, so running just once would most likely not result in an accurate estimation.

3. Assuming we are trying to generate one cluster for each Iris species represented in the dataset, does one measure seem to be a consistently better criterion for choosing the value of k than the other? Why or why not?

The SSQ elbow measurement seems to be consistently better for choosing the value of k than gap statistics. However, gap statistics provide more data, but this can make it more difficult to choose the value of k without further data processing of the statistical output. By visual comparison of the "Clustering Sums of Squared Distances" and "Clustering Gap Statistics" we can see that the elbow is in a similar position, but the tailing data in gap statistics can potentially be misinterpreted. This dataset seems to work well for both.

