

In [1]:

```
1 from nltk.tokenize import word_tokenize
2 from nltk.corpus import stopwords
3 from string import punctuation
4 import re
```

In [2]:

```
1 documents = [
2     "This little Kitty came to play when I was eating at a restaurant.",
3     "Merley has the best squooshy kitten belly.",
4     "Google Translate app is incredible",
5     "If you open 100 tab in google you get a smiley face.",
6     "Best cat photo I've ever taken.",
7     "Climbing ninja cat.",
8     "Impressed with google map feedback",
9     "Key promoter extension for Google Chrome."
10 ]
```

In [3]:

```
1 from collections import defaultdict
2
3 DTM = defaultdict(lambda:defaultdict(int))
4 for i,d in enumerate(documents):
5     for t in word_tokenize(d.lower()):
6         if t not in stopwords.words("english") and not re.search(r"%s" % re.escape(punctuation), t):
7             DTM[i][t] += 1
```

In [6]:

```
1 TDM = defaultdict(lambda:defaultdict(int))
2 for d, termList in DTM.items():
3     for t, f in termList.items():
4         TDM[t][d] = f
```

In []:

```
1 '''
2 TF * IDF
3 Tf = freq(i,j) / max freq(j) => j번째 문서에서 나온, i번째 빈도를 freq(i,j)
4 freq(i,j) => TDM[단어][문서:빈도], max freq => DTM[문서] max([단어:빈도])
5 IDF = log N / df(i) => i번째 단어의 문서빈도 df(i), N => 전체 문서 수
6 df => len(TDM[단어]), N => len(DTM)
7 '''
```

In [8]:

```
1 from math import log2
2
3 TWM = defaultdict(lambda:defaultdict(float))
4 N = len(DTM)
5 for term, docList in TDM.items():
6     df = len(docList)
7     for d,f in docList.items():
8         maxFreq = max(DTM[d].values())
9         TF = f/maxFreq
10        IDF = log2(N/df)
11        TWM[term][d] = TF*IDF
12 #TWM[단어][문서:가중치,...]
```

In [9]:

1	TWM
---	-----

...

K-Means

In []:

```
1 '''
2 TWM => 단어-문서:가중치, ...
3 단어 -> concept(차원 축)
4 문서->BoW(Bag of words)[단어,단어,단어, .. V=TDM.keys()]
5 문서1 = [0,0,1,0,0,..., V] <- vectorize
6 문서2 = [1,0,1,0,0,..., V]
7
8 '''
9 K = ?
10 centroid = ? <- k개 만큼.. [for _ in range(K)]
11 _iter = 10
12 for _ in range(_iter):
13     expectation()
14     maximization()
15 '''
```

In [12]:

```
1 docVectorList = [[0 for _ in range(len(TDM))] for _ in range(N)]
2 D = list(DTM.keys())
3 V = list(TDM.keys())
4 for t, docList in TWM.items():
5     for d,w in docList.items():
6         docVectorList[D.index(d)][V.index(t)] = w
```

In [71]:

```
1 | V
```

Out[71]:

```
['little',
 'kitty',
 'came',
 'play',
 'eating',
 'restaurant',
 'merley',
 'best',
 'squoooshy',
 'kitten',
 'belly',
 'google',
 'translate',
 'app',
 'incredible',
 'open',
 '100',
 'tab',
 'get',
 'smiley',
 'face',
 'cat',
 'photo',
 'ever',
 'taken',
 'climbing',
 'ninja',
 'impressed',
 'map',
 'feedback',
 'key',
 'promoter',
 'extension',
 'chrome']
```

In [16]:

```
1 | print(docVectorList[0])
2 | print(D[0], V[:6])
```

```
[3.0, 3.0, 3.0, 3.0, 3.0, 3.0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0]
0 ['little', 'kitty', 'came', 'play', 'eating', 'restaurant']
```

In [11]:

```
1 | len(docVectorList), len(documents), len(docVectorList[0]), len(TDM)
```

Out[11]:

```
(8, 8, 34, 34)
```

In [52]:

```
1 | from math import sqrt
2 |
3 | def distance(x1, x2):
4 |     _sum = 0
5 |     for i in range(len(x1)):
6 |         _sum += (x1[i]-x2[i])**2
7 |     return sqrt(_sum)
8 |
9 | def angle(x1, x2):
10 |    _innerProduct = 0
11 |    for i in range(len(x1)):
12 |        _innerProduct += x1[i]*x2[i]
13 |    x1VecLength = distance(x1, [0 for _ in range(len(x1))])
14 |    x2VecLength = distance(x2, [0 for _ in range(len(x2))])
15 |    return _innerProduct/(x1VecLength*x2VecLength)
```

In [51]:

```
1 | def expectation(x, c, opt=False):
2 |     candidates = list()
3 |
4 |     nearest = distance if not opt else angle
5 |     best = min if not opt else max
6 |
7 |     for _ in c:
8 |         candidates.append(nearest(x, _))
9 |     return candidates.index(best(candidates))
```

In [47]:

```
1 | def maximization(X):
2 |     _sum = [0 for _ in range(len(X[0]))]
3 |     D = len(X)
4 |     for x in X:
5 |         for i in range(len(x)):
6 |             _sum[i] += x[i]
7 |     return [_/D for _ in _sum]
```

In [49]:

```
1 | def sse(X, c, opt=False): # sum sward error
2 |     error = 0
3 |     nearest = distance if not opt else angle
4 |     for x in X:
5 |         error += nearest(x,c)
6 |     return error
```

In [54]:

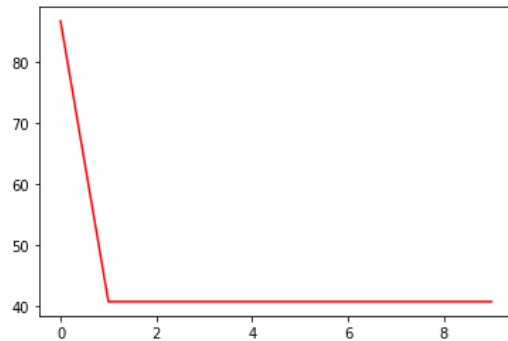
```
1 from random import randrange
2
3 K = 2
4 centroid = [[randrange(0,4) for _ in range(len(V))] for _ in range(K)]
5 _iter = 10
6 sseList = list()
7
8 for _ in range(_iter):
9     rnk = list(list(0 for _ in range(K)) for _ in range(len(D)))
10    for i,x in enumerate(docVectorList):
11        idx = expectation(x, centroid)
12        rnk[i][idx] = 1
13
14    _sum = 0.0
15    for k in range(K):
16        _X = [docVectorList[i] for i in range(len(D)) if rnk[i][k]]
17        _sum += sse(_X, centroid[k]) #sum squared error
18        centroid[k] = maximization(_X)
19    sseList.append(_sum)
```

In [57]:

```
1 import matplotlib.pyplot as plt
2
3 plt.plot(range(_iter), sseList, "r-")
4 plt.show
```

Out[57]:

<function matplotlib.pyplot.show(*args, **kw)>



In [58]:

```
1 for k in range(K):
2     _X = [i for i in range(len(D)) if rnk[i][k]]
3     print("Cluster:",k)
4     for x in _X:
5         print(documents[x])
```

Cluster: 0

This little kitty came to play when I was eating at a restaurant.

If you open 100 tab in google you get a smiley face.

Cluster: 1

Merley has the best squooshy kitten belly.

Google Translate app is incredible

Best cat photo I've ever taken.

Climbing ninja cat.

Impressed with google map feedback

Key promoter extension for Google Chrome.

In [59]:

```
1 [{V[i]:centroid[0][i]} for i in range(len(V))]
```

...

In [61]:

```
1 for k in range(K):
2     print([V[i]:centroid[k][i]} for i in range(len(V)) if centroid[k][i]])
```

```
[{'little': 1.5}, {'kitty': 1.5}, {'came': 1.5}, {'play': 1.5}, {'eating': 1.5}, {'restaurant': 1.5}, {'google': 0.5}, {'open': 1.5}, {'100': 1.5}, {'tab': 1.5}, {'get': 1.5}, {'smiley': 1.5}, {'face': 1.5}]
[{'merley': 0.5}, {'best': 0.6666666666666666}, {'squooshy': 0.5}, {'kitten': 0.5}, {'belly': 0.5}, {'google': 0.5}, {'translate': 0.5}, {'app': 0.5}, {'incredible': 0.5}, {'cat': 0.6666666666666666}, {'photo': 0.5}, {'ever': 0.5}, {'taken': 0.5}, {'climbing': 0.5}, {'ninja': 0.5}, {'impressed': 0.5}, {'map': 0.5}, {'feedback': 0.5}, {'key': 0.5}, {'promoter': 0.5}, {'extension': 0.5}, {'chrome': 0.5}]
```

In [65]:

```
1 for k in range(K):
2     W = [(V[i].centroid[k][i]) for i in range(len(V)) if centroid[k][i]]
3     print(sorted(W, key=lambda x:x[1], reverse=True)[:5])
```

```
[('little', 1.5), ('kitty', 1.5), ('came', 1.5), ('play', 1.5), ('eating', 1.5)]
[('best', 0.6666666666666666), ('cat', 0.6666666666666666), ('merley', 0.5), ('squooshy', 0.5), ('kitten', 0.5)]
```

In [69]:

```
1 result = list()
2 for k in range(K):
3     W = [(V[i],centroid[k][i]) for i in range(len(V)) if centroid[k][i]]
4     result.append(sorted(W, key=lambda x:x[1], reverse=True)[:5])
5     #print(sorted(W, key=lambda x:x[1], reverse=True)[:5])
6 result
```

Out[69]:

```
[(['little', 1.5),
 ('kitty', 1.5),
 ('came', 1.5),
 ('play', 1.5),
 ('eating', 1.5)],
 [('best', 0.6666666666666666),
 ('cat', 0.6666666666666666),
 ('merley', 0.5),
 ('squoochy', 0.5),
 ('kitten', 0.5)]]
```

In [70]:

```
1 from wordcloud import WordCloud
2
3 data = {x[0]:x[1] for x in result[0]}
4 wc = WordCloud(background_color="white")
5 wc.generate_from_frequencies(data)
6 wc.to_image()
```

Out[70]:



In []:

```
1 V = len(distinct words:Vocaburary)
2 K = 6
3 from random import sample
4 len(sample(docVeectorlist, K))
```