# 读取数据，图片类型

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
from sklearn.metrics import pairwise_distances_argmin
from sklearn.datasets import load_sample_image
from sklearn.utils import shuffle
from time import time
from numpy import *
import time
import matplotlib.pyplot as plt

n_colors = 64
n_clusters=n_colors
# Load the Summer Palace photo
china = load_sample_image("china.jpg")
print(china.shape)

china = np.array(china, dtype=np.float64) / 255
w, h, d = original_shape = tuple(china.shape)
image_array = np.reshape(china, (w * h, d))
image_array_sample = shuffle(image_array, random_state=0)[:1000]
```

```
(427, 640, 3)
```

```
C:\Users\Magneto_Wang\Anaconda3\lib\site-packages\sklearn\datasets\base.py:762:
DeprecationWarning: `imread` is deprecated!
`imread` is deprecated in SciPy 1.0.0, and will be removed in 1.2.0.
Use ``imageio.imread`` instead.
  images = [imread(filename) for filename in filenames]
C:\Users\Magneto_Wang\Anaconda3\lib\site-packages\sklearn\datasets\base.py:762:
DeprecationWarning: `imread` is deprecated!
`imread` is deprecated in SciPy 1.0.0, and will be removed in 1.2.0.
Use ``imageio.imread`` instead.
  images = [imread(filename) for filename in filenames]
```

```
plt.figure(4)
plt.clf()
ax = plt.axes([0, 0, 1, 1])
plt.axis('off')
plt.title('原图')
plt.imshow(china)
```

```
1  <matplotlib.image.AxesImage at 0x20b8fff59b0>
```



# 基本函数的构造

# kmeans

```python
1   # k-means cluster
2   #dataSet为一个矩阵
3   #k为将dataSet矩阵中的样本分成k个类
4   def kmeans(dataSet, k):
5       numSamples = dataSet.shape[0]   #读取矩阵dataSet的第一维度的长度,即获得有多少个样本数据
6       # first column stores which cluster this sample belongs to,
7       # second column stores the error between this sample and its centroid
8       clusterAssment = mat(zeros((numSamples, 2)))   #得到一个N*2的零矩阵
9       clusterChanged = True
10
11      ## step 1: init centroids
12      centroids = initCentroids(dataSet, k)   #在样本集中随机选取k个样本点作为初始质心
13
14      while clusterChanged:
15          clusterChanged = False
16          ## for each sample
17          for i in range(numSamples):   #range
18              minDist  = 100000.0
19              minIndex = 0
20                  ## for each centroid
```

```
21            ## step 2: find the centroid who is closest
22            #计算每个样本点与质点之间的距离，将其归内到距离最小的那一簇
23            for j in range(k):
24                distance = euclDistance(centroids[j, :], dataSet[i, :])
25                if distance < minDist:
26                    minDist  = distance
27                    minIndex = j
28
29            ## step 3: update its cluster
30            #k个簇里面与第i个样本距离最小的的标号和距离保存在clusterAssment中
31            #若所有的样本不在变化，则退出while循环
32            if clusterAssment[i, 0] != minIndex:
33                clusterChanged = True
34                clusterAssment[i, :] = minIndex, minDist**2   #两个**表示的是minDist的平方
35
36        ## step 4: update centroids
37        for j in range(k):
38            #clusterAssment[:,0].A==j是找出矩阵clusterAssment中第一列元素中等于j的行的下标，返回
    的是一个以array的列表，第一个array为等于j的下标
39            pointsInCluster = dataSet[nonzero(clusterAssment[:, 0].A == j)[0]] #将dataSet矩阵
    中相对应的样本提取出来
40            centroids[j, :] = mean(pointsInCluster, axis = 0)   #计算标注为j的所有样本的平均值
41
42    print ('Congratulations, cluster complete!')
43    return centroids, clusterAssment
44
45 # calculate Euclidean distance
46 def euclDistance(vector1, vector2):
47    return sqrt(sum(power(vector2 - vector1, 2)))   #求这两个矩阵的距离，vector1、2均为矩阵
48
49 def initCentroids(dataSet, k):
50    numSamples, dim = dataSet.shape    #矩阵的行数、列数
51    centroids = zeros((k, dim))          #感觉要不要你都可以
52    for i in range(k):
53        index = int(random.uniform(0, numSamples))  #随机产生一个浮点数，然后将其转化为int型
54        centroids[i, :] = dataSet[index, :]
55    return centroids
56
```

# 原图各个像素的分类

```
1 def predict_data(dataSet, k,centroids, clusterAssment):
2    numSamples = dataSet.shape[0]
3    label=[]
4    predict_Assment = mat(zeros((numSamples, 2)))
5
6    clusterChanged = True
7    for i in range(numSamples):
8        minDist  = 100000.0
9        minIndex = 0
10        for j in range(k):
```

```
11                distance = euclDistance(centroids[j, :], dataSet[i, :])
12            if distance < minDist:
13                minDist  = distance
14                minIndex = j
15            if predict_Assment[i, 0] != minIndex:
16                clusterChanged = True
17                predict_Assment[i, :] = minIndex, minDist**2
18
19
20        return predict_Assment
```

# 重构图片

```
1  def recreate_image(codebook, labels, w, h):
2      """Recreate the (compressed) image from the code book & labels"""
3      d = codebook.shape[1]
4      image = np.zeros((w, h, d))
5      label_idx = 0
6      for i in range(w):
7          for j in range(h):
8              #labelIndex,dist= labels[i, :]
9              image[i][j] = codebook[int(label[label_idx, :][0,:1])]
10             label_idx += 1
11     return image
12
```

# 图片数据量非常大。为了更快计算。我是抽样其中部分数据。虽然图像有损失，但是效果依然不错

## 聚类的点数是自定义。这里选择3，30，60

```
1  centroids, clusterAssment=kmeans(image_array_sample, 3)
```

```
1  label=predict_data(image_array,3,centroids, clusterAssment)
2  image=recreate_image(centroids,label,w,h)
```
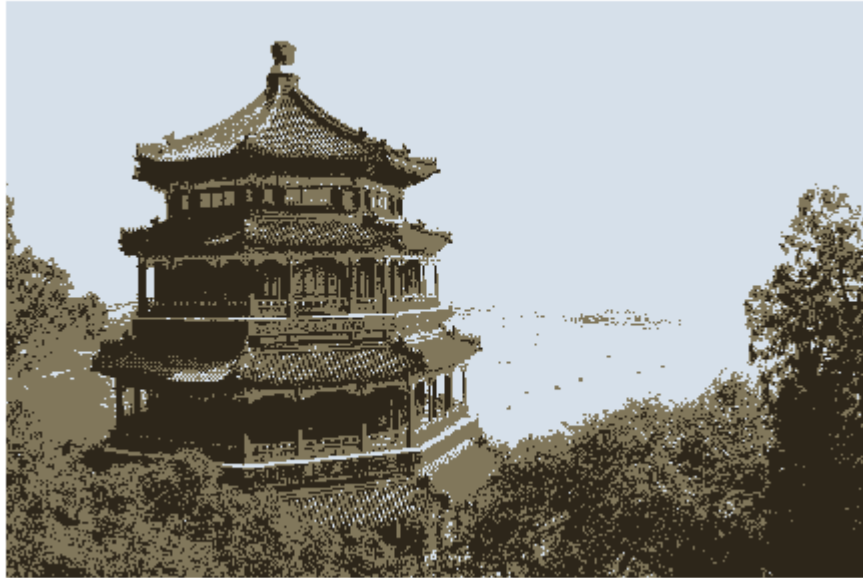
```
1  plt.figure(1)
2  plt.clf()
3  ax = plt.axes([0, 0, 1, 1])
4  plt.axis('off')
5  plt.title('Quantized image (3 colors, K-Means)')
6  plt.imshow(image)
```

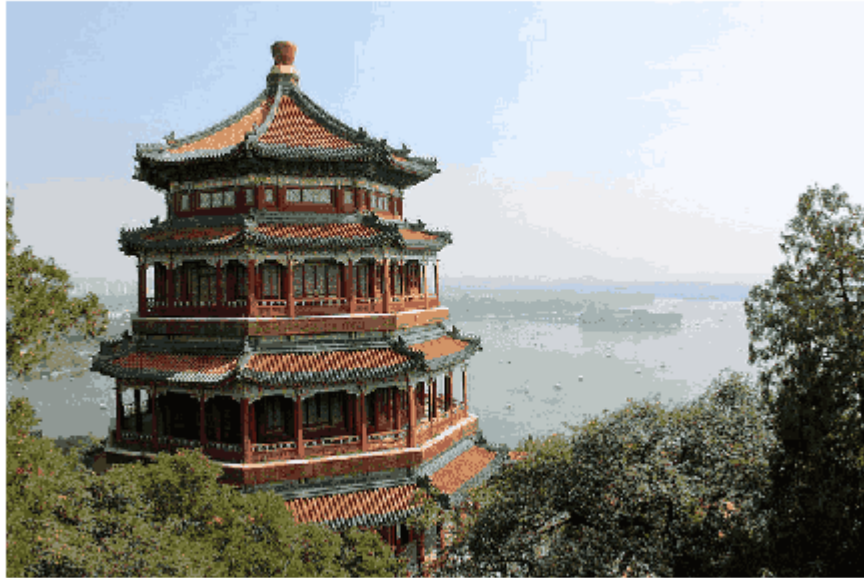## Quantized image (3 colors, K-Means)



```
1  centroids, clusterAssment=kmeans(image_array_sample, 30)
```

```
1  label=predict_data(image_array,30,centroids, clusterAssment)
2  image=recreate_image(centroids,label,w,h)
```

```
1  plt.figure(2)
2  plt.clf()
3  ax = plt.axes([0, 0, 1, 1])
4  plt.axis('off')
5  plt.title('Quantized image (30 colors, K-Means)')
6  plt.imshow(image)
```

Quantized image (30 colors, K-Means)



```
1   n_color=60
2   centroids, clusterAssment=kmeans(image_array_sample, n_color)
3   label=predict_data(image_array,n_color,centroids, clusterAssment)
4   image=recreate_image(centroids,label,w,h)
5   plt.figure(3)
6   plt.clf()
7   ax = plt.axes([0, 0, 1, 1])
8   plt.axis('off')
9   plt.title('Quantized image (60 colors, K-Means)')
10  plt.imshow(image)
```

```
1   Congratulations, cluster complete!
```

Quantized image (60 colors, K-Means)