

k-Means, PCA, and Tree VQ

Thomas Breuel

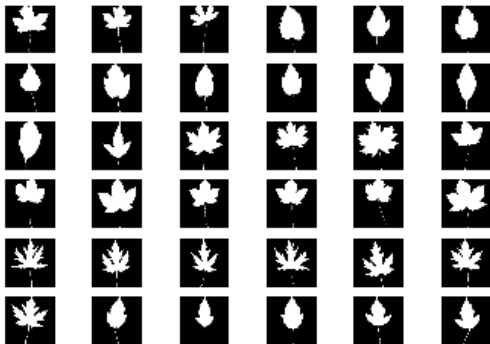
UniKL

```
1 import random as pyrand
2 import tables
3 matplotlib.rc("image", cmap="gray")
```

Clustering Shapes

```
1 hdf = tables.openFile("leaves.h5","r")
2 nimages = len(hdf.root.icons)
3 images = hdf.root.icons[:, :, :] / 256.0
4 data = images.reshape(nimages, -1)
```

```
1 for i in range(36):  
2     subplot(6,6,i+1); imshow(images[i]); axis("off")
```

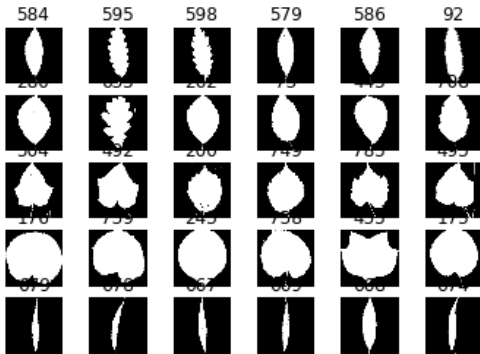


```
1 from sklearn.cluster import KMeans, MiniBatchKMeans
2 km = MiniBatchKMeans(n_clusters=36)
3 km.fit(data)
4 assignments = km.predict(data)
5 nearest = [[i for i, j in enumerate(assignments) if j==c] for c in
             range(36)]
6 good = [1 for l in nearest if len(l)>3]
```

```

1 for i,l in enumerate(good[:5]):
2     for j,k in enumerate(pyrand.sample(1,min(6,len(l)))):
3         subplot(5,6,i*6+j+1)
4         imshow(hdf.root.icons[k])
5         title("%d"%(k,))
6         axis("off")

```



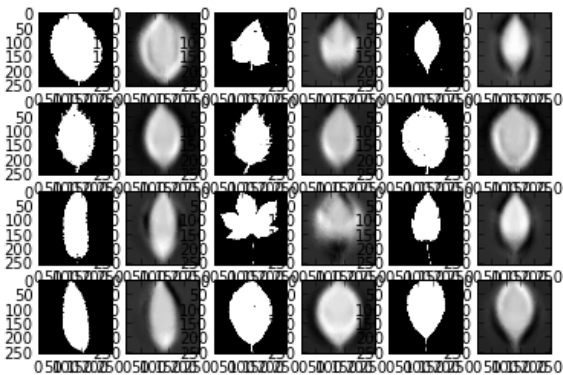
Combining PCA and k -means


```
1 from sklearn.decomposition import RandomizedPCA
2 pca = RandomizedPCA(n_components=10)
3 pca.fit(data)
4 lo = pca.transform(data)
5 hi = pca.inverse_transform(lo)
```

```

1 for i,j in enumerate(pyrand.sample(range(len(images)),12)):
2     subplot(4,6,2*i+1); imshow(images[j])
3     subplot(4,6,2*i+2); imshow(hi[j].reshape(256,256))

```

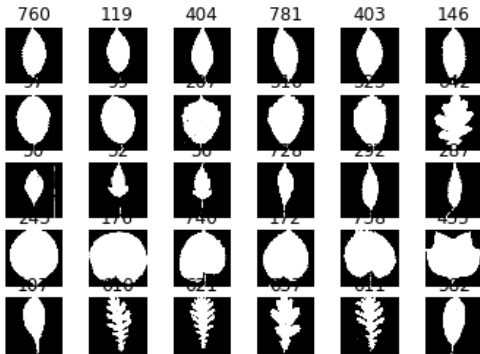


```
1 from sklearn.cluster import KMeans, MiniBatchKMeans
2 km = MiniBatchKMeans(n_clusters=36)
3 km.fit(lo)
4 assignments = km.predict(lo)
5 nearest = [[i for i,j in enumerate(assignments) if j==c] for c in
6            range(36)]
7 good = [1 for l in nearest if len(l)>3]
```

```

1 for i,l in enumerate(good[:5]):
2     for j,k in enumerate(pyrand.sample(1,min(6,len(l)))):
3         subplot(5,6,i*6+j+1)
4         imshow(images[k])
5         title("%d"%(k,))
6         axis("off")

```



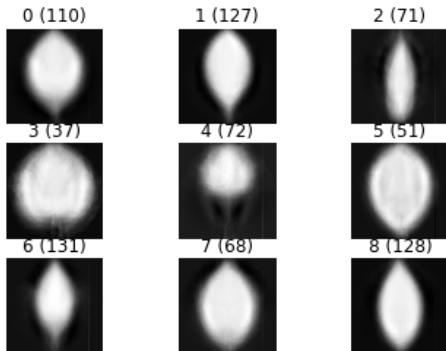
Hierarchical and Tree VQ

```
1 km = MiniBatchKMeans(n_clusters=9)
2 km.fit(lo)
3 assignments = km.predict(lo)
4 nearest = [[i for i,j in enumerate(assignments) if j==c] for c in
             range(9)]
```

```

1 for i,c in enumerate(km.cluster_centers_):
2     im = pca.inverse_transform(c)
3     subplot(3,3,i+1); axis("off"); imshow(im.reshape(256,256))
4     title("%d_␣(%d)"%(i,len(nearest[i])))

```



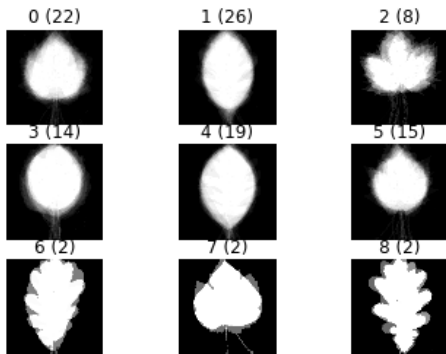
```
1 data0 = [data[i] for i in nearest[0]] # note data instead of lo
2 print len(data0)
3 km0 = MiniBatchKMeans(n_clusters=9)
4 km0.fit(data0)
5 assignments0 = km0.predict(data0)
6 nearest0 = [[i for i,j in enumerate(assignments0) if j==c] for c in
              range(9)]
```

110


```

1 for i,c in enumerate(km0.cluster_centers_):
2     subplot(3,3,i+1); axis("off"); imshow(c.reshape(256,256))
3     title("%d_␣(%d)"%(i,len(nearest0[i])))

```



Hierarchical Tree VQ

- ▶ we perform VQ in several stages
- ▶ a first, top-level VQ
- ▶ then we split the data based on the top-level assignments
- ▶ then we repeat the VQ for each split
- ▶ this can be continued recursively (tree VQ)
- ▶ if we use different resolutions or number of PCA components at each level, we are using a hierarchical VQ

Hierarchical tree VQ is a good way of clustering very large data sets.