

ONLINE MASTERS IN DATA SCIENCE

DSC 257R - UNSUPERVISED LEARNING

TWO USES OF CLUSTERING

SANJOY DASGUPTA, PROFESSOR

UC San Diego

COMPUTER SCIENCE & ENGINEERING
HALICIOĞLU DATA SCIENCE INSTITUTE

Two Common Uses of Clustering

- Finding meaningful structure in data
Finding salient grouping in data.
- Vector quantization
Find a finite set of representatives that provides good coverage of a complex, possibly infinite, high-dimensional space.

Representing Images Using K -Means Codewords

How to represent a collection of images as fixed-length vectors?



Representing Images Using K -Means Codewords

How to represent a collection of images as fixed-length vectors?



- Take all $\ell \times \ell$ patches in all images. Extract features for each.
- Run k -means on this entire collection to get k centers.
- Now associate any image patch with its nearest center.
- Represent an image by a histogram over $\{1, 2, \dots, k\}$.

Looking for Natural Groups in Data

"Animals with attributes" data set

- 50 animals: antelope, grizzly bear, beaver, dalmatian, tiger, ...
- 85 attributes: longneck, tail, walks, swims, nocturnal, forager, desert, bush, plains, ...
- Each animal gets a score (0 - 100) along each attribute
- 50 data points in \mathbb{R}^{85}

Apply k -means with $k = 10$ and look at grouping obtained.

- 1 zebra
- 2 spider monkey, gorilla, chimpanzee
- 3 tiger, leopard, wolf, bobcat, lion
- 4 hippopotamus, elephant, rhinoceros
- 5 killer whale, blue whale, humpback whale, seal, walrus, dolphin
- 6 giant panda
- 7 skunk, mole, hamster, squirrel, rabbit, bat, rat, weasel, mouse, raccoon
- 8 antelope, horse, moose, ox, sheep, giraffe, buffalo, deer, pig, cow
- 9 beaver, otter
- 10 grizzly bear, dalmatian, persian cat, german shepherd, siamese cat, fox, chihuahua, polar bear, collie

- 1 zebra
- 2 spider monkey, gorilla, chimpanzee
- 3 tiger, leopard, fox, wolf, bobcat, lion
- 4 hippopotamus, elephant, rhinoceros, buffalo, pig
- 5 killer whale, blue whale, humpback whale, seal, otter, walrus, dolphin
- 6 dalmatian, persian cat, german shepherd, siamese cat, chihuahua, giant panda, collie
- 7 beaver, skunk, mole, squirrel, bat, rat, weasel, mouse, raccoon
- 8 antelope, horse, moose, ox, sheep, giraffe, deer, cow
- 9 hamster, rabbit
- 10 grizzly bear, polar bear