

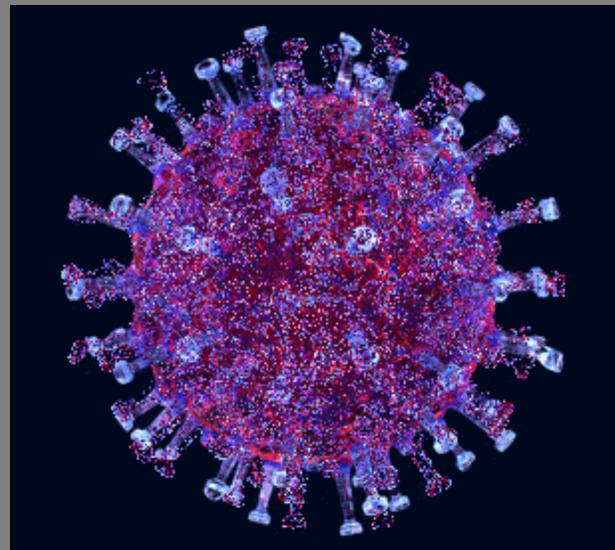
COS324: Introduction to Machine Learning

Topic: Unsupervised Learning, Clustering, Embeddings

Prof. Yoram Singer

April 28, 2020

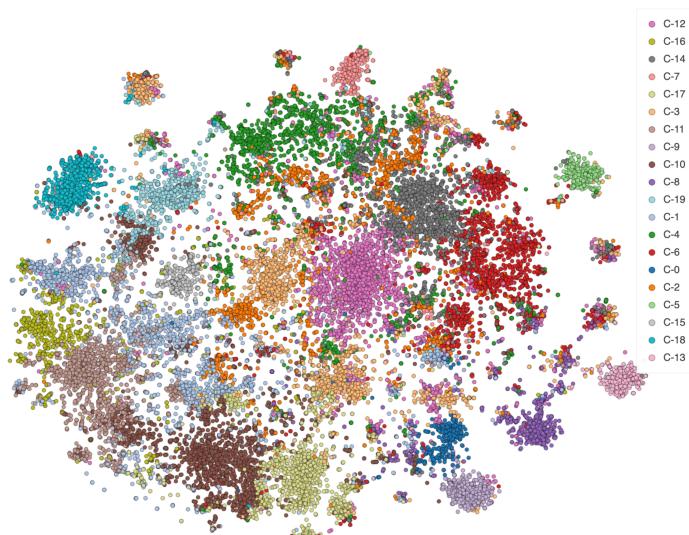
Pressing Problem



1 / 32

2 / 32

Clustering of Scientific Articles on Covid-19



3 / 32

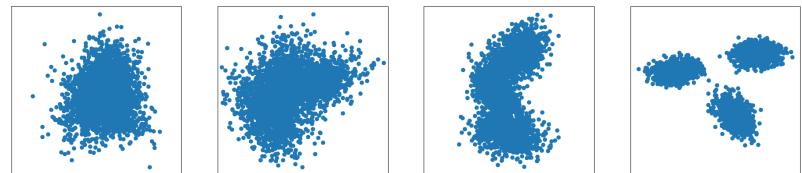
Unsupervised Learning

- So far discussed supervised learning:
 - Examples (\mathbf{x}, y) are input-target pairs in $\mathcal{X} \times \mathcal{Y}$
 - Learning amounts to finding a function $f : \mathcal{X} \rightarrow \mathcal{Y}$
 - Loss measures discrepancy between y and $\hat{y} = f(\mathbf{x})$: $\ell(y, \hat{y})$
- Sometimes we have plentiful of instances \mathbf{x}_i
 - ... but only handful of labels $m \gg n$
 - ... or none at all $n = 0$
- It is nonetheless useful to find “structure” or meaningful patterns in the data

4 / 32

Goals of Clustering

- Intuitively, grouping a set of objects such that
 - Similar objects end up in the same cluster
 - Dissimilar objects place in different clusters
- Imprecise & potentially ambiguous definition
- Disappointingly, no single simple definition



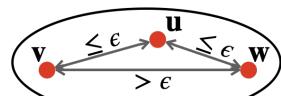
5 / 32

6 / 32

Sources of Difficulty

- Inherit problem: lack of “ground truth” and tangible objective
- Technical difficulty:
 - Similarity & distance functions are not transitive

$$\|\mathbf{u} - \mathbf{v}\| \leq \epsilon \wedge \|\mathbf{u} - \mathbf{w}\| \leq \epsilon \not\Rightarrow \|\mathbf{v} - \mathbf{w}\| \leq \epsilon$$

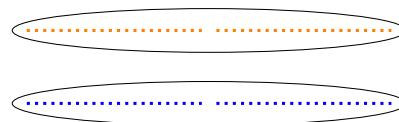


- Cluster (set) membership is transitive
 - Define $\mathbf{u} \sim \mathbf{v}$ iff \mathbf{u} and \mathbf{v} belong to the same cluster
 - Then, $\mathbf{u} \sim \mathbf{v} \wedge \mathbf{v} \sim \mathbf{w} \Rightarrow \mathbf{u} \sim \mathbf{w}$

Clustering is Ambiguous

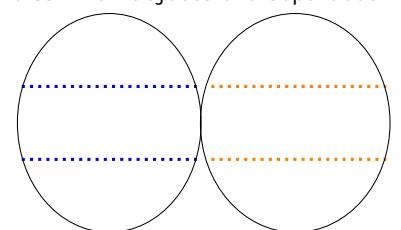
.....

similar objects in same cluster



.....

dissimilar objects are separated

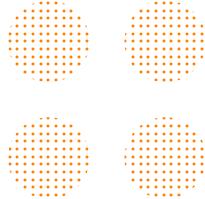


7 / 32

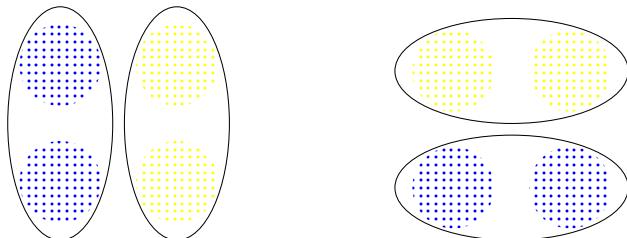
8 / 32

Lack of Ground Truth

Partition points into **two** clusters:



We have two well justifiable solutions:



9 / 32

10 / 32

Clustering Problem Setting

- Input: $S = \{\mathbf{x}_i\}_{i=1}^m$ where $\mathbf{x}_i \in \mathbb{R}^d$
- Distance $d : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}_+$ or similarity $s : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ [s need not be symmetric, $s(\mathbf{u}, \mathbf{v}) \neq s(\mathbf{v}, \mathbf{u})$]
- Output: partition $\mathcal{C} = \{C_i\}_{i=1}^k$ of set S into k subsets s.t.

$$S = \bigcup_{i=1}^k C_i \quad C_i \cap C_j = \emptyset$$

- Number of clusters k may be part of input or unknown

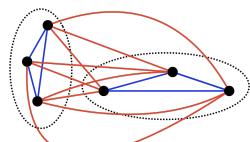
Cost-based Clustering

- We focus on distance-based $d(\mathbf{u}, \mathbf{v})$ clustering
- Cost of partitioning $\mathcal{C} = \{C_i\}_{i=1}^k$ of S ?
- Define indicator

$$\mathbb{1}[i, j | \mathcal{C}] = \begin{cases} +1 & \exists r : \mathbf{x}_i \in C_r \wedge \mathbf{x}_j \in C_r \\ -1 & \text{o.w.} \end{cases}$$

- Penalize for large intra-cluster & small inter-cluster distances

$$\ell(S, \mathcal{C}) = \sum_{i, j=1}^{|S|} \mathbb{1}[i, j | \mathcal{C}] d(\mathbf{x}_i, \mathbf{x}_j)$$



- Number of instances to compare $O(n^2)$

$$|C_i| \approx \frac{n}{k} \Rightarrow \binom{k}{2} \left(\frac{n}{k} \right)^2 \equiv O(n^2) \quad \text{inter-cluster pairs}$$

$$k \binom{n}{2} \equiv O\left(\frac{n^2}{k}\right) \quad \text{intra-cluster pairs}$$

11 / 32

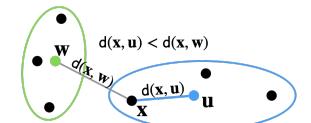
k-Center Clustering

- Centroid-based clustering: intuitive, transitive, "aesthetic"
- Associate a center $\mathbf{w}_j \in \mathbb{R}^d$ with partition C_j

$$\mathbf{x}_i \in C_j \Leftrightarrow \forall r \neq j : d(\mathbf{x}_i, \mathbf{w}_j) < d(\mathbf{x}_i, \mathbf{w}_r)$$

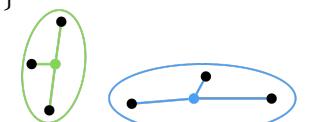
- Induces partitioning

$$C_j = \{i : \forall r \neq j d(\mathbf{x}_i, \mathbf{w}_j) < d(\mathbf{x}_i, \mathbf{w}_r)\}$$



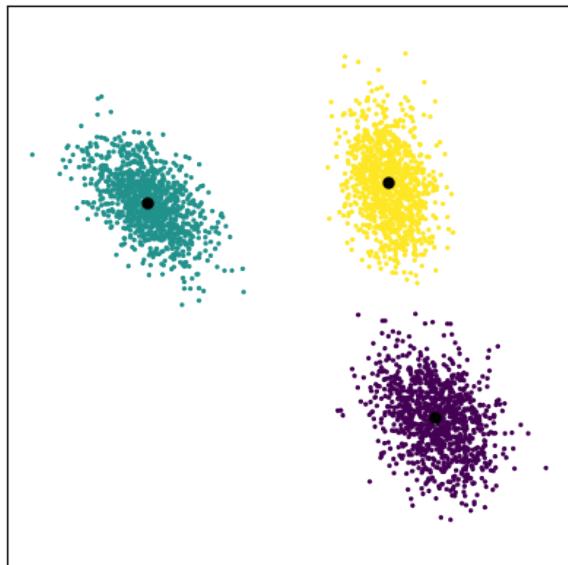
- Loss of k-centers

$$\ell(S, \mathcal{C}) = \sum_{j=1}^k \sum_{i \in C_j} d(\mathbf{x}_i, \mathbf{w}_j) = \sum_{i=1}^m \min_{j=1}^k d(\mathbf{x}_i, \mathbf{w}_j)$$



12 / 32

Example of 3-Center Clustering



13 / 32

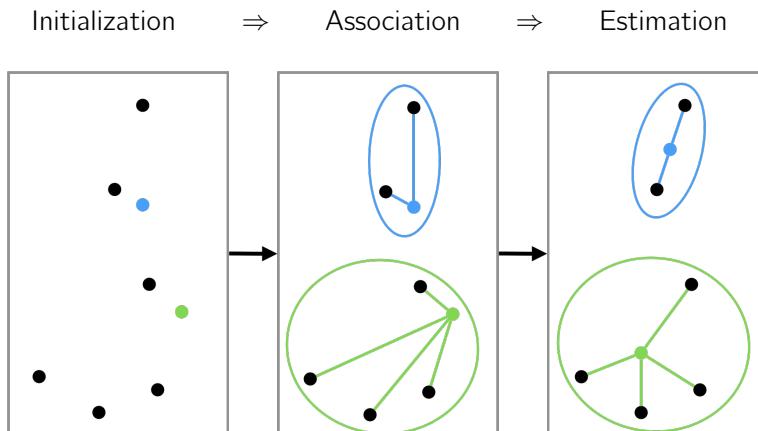
Skeleton of Metric Clustering

- Initialize each \mathbf{w}_j^0 to a vector in \mathbb{R}^d
 - For $t = 1, \dots, T$:
 - Associate each \mathbf{x}_i with its nearest centroid
$$\forall i : a^t(i) = \arg \min_{j=1}^k d(\mathbf{x}_i, \mathbf{w}_j^{t-1})$$
 - Restimate centroids from associations
- $$\forall j : \mathbf{w}_j^t = \min_{\mathbf{w}} \sum_{i: a^t(i)=j} d(\mathbf{x}_i, \mathbf{w})$$

- If $\forall i : a^t(i) = a^{t-1}(i)$ break

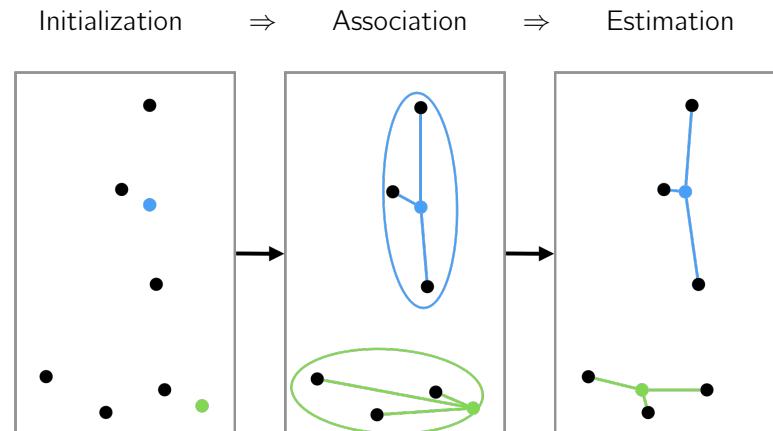
14 / 32

Initialization makes a difference: Example I.a



15 / 32

Initialization makes a difference: Example I.b

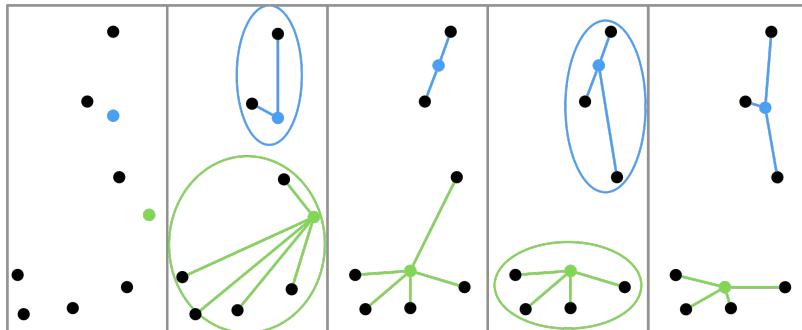


16 / 32

Example II

A single example might change clustering entirely:

Initialization \Rightarrow Associate \Rightarrow Estimate \Rightarrow Associate \Rightarrow Estimate



17 / 32

Convergence of Metric-based Clustering

- Centers at iteration t

$$\mathcal{W}^t = \{\mathbf{w}_j^t\}_{j=1}^k$$

- Partition at iteration t

$$\mathcal{A}^t = \{a^t(i)\}_{i=1}^m$$

- Loss of partition and centers

$$\ell(S, \mathcal{A}, \mathcal{W}) = \frac{1}{m} \sum_{i=1}^m d(\mathbf{x}_i, \mathbf{w}_{a(i)})$$

- Then, $\ell(S, \mathcal{A}^{t-1}, \mathcal{W}^{t-1}) > \ell(S, \mathcal{A}^t, \mathcal{W}^{t-1}) > \ell(S, \mathcal{A}^t, \mathcal{W}^t)$
- Since $\ell(S, \mathcal{A}, \mathcal{W}) \geq 0$ and $\forall t, j : \mathbf{w}_j^t \in \bar{S}$
 $\Rightarrow \ell(S, \mathcal{A}^t, \mathcal{W}^t)$ converges to a local minimum

18 / 32

k-Means

- Use $d(\mathbf{u}, \mathbf{v}) \stackrel{\text{def}}{=} \|\mathbf{u} - \mathbf{v}\|_2^2$
- Solving $\min_{\mathbf{w}} \sum_{i:a(i)=j} \|\mathbf{x}_i - \mathbf{w}\|^2$ amounts to

$$\mathbf{w}_j = \frac{1}{n_j} \sum_{i:a(i)=j} \mathbf{x}_i \quad \text{where } n_j \stackrel{\text{def}}{=} |\{i : a(i) = j\}|$$

- Namely, center of mass of examples in cluster
- Runtime is: Tkn

k-Medians

- Use $d(\mathbf{u}, \mathbf{v}) \stackrel{\text{def}}{=} \|\mathbf{u} - \mathbf{v}\|_1$
- Solving $\min_{\mathbf{w}} \sum_{i:a(i)=j} \|\mathbf{x}_i - \mathbf{w}\|_1$ amounts to

$$\mathbf{w}_j[r] = \min_{\omega} \sum_{i:a(i)=j} |\mathbf{x}_i[r] - \omega| \\ = \text{median}\{\mathbf{x}_i[r] : a(i) = j\}$$
- $\mathbf{w}_j[r]$ is median of r 'th coordinate of examples in cluster
- Runtime is: Tkn

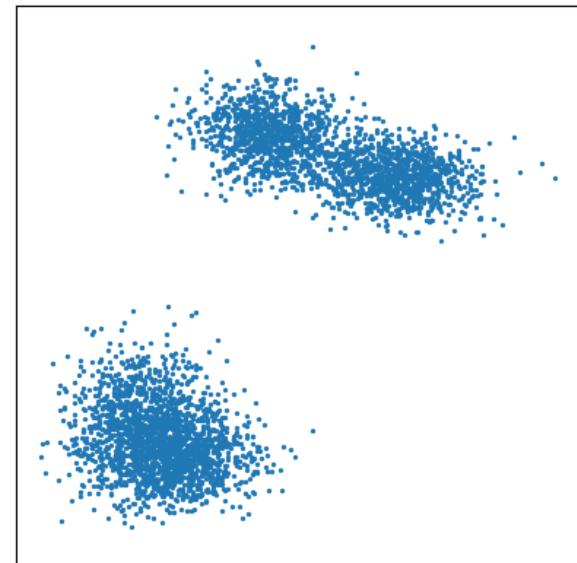
19 / 32

20 / 32

Tricks & Treats

- Initialization:
 - At random
 - Agglomeratively: warm-start from $k - 1$ clusters
 - Agglomeratively: hierarchical from $2 \times \frac{k}{2}$ clusters
 - Using other clustering methods (e.g. spectral)
- Art of choosing number of clusters k ...
- Small amounts of labeled data:
 - Determine number of clusters
 - Good initialization
 - Metric adjustment prior to clustering

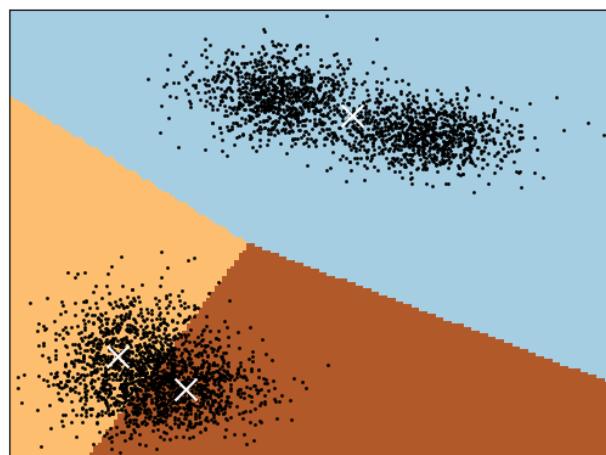
Data Generated by k Gaussians



21 / 32

22 / 32

Clustering with $\hat{k} = 3$

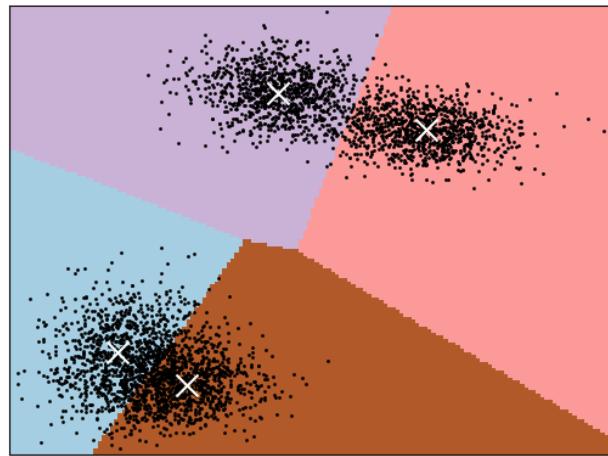


Why do we see straight decision boundaries ?

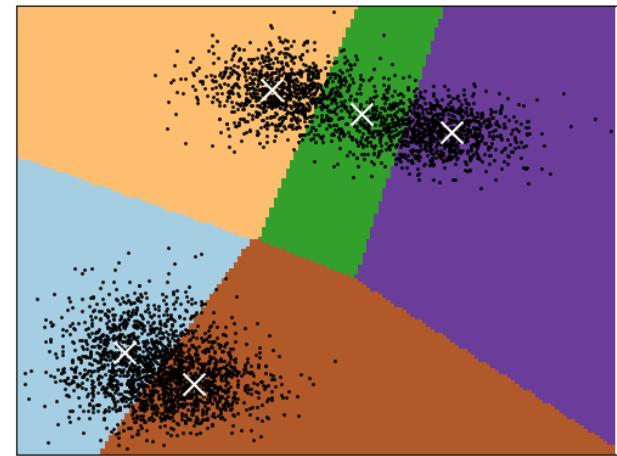
23 / 32

24 / 32

Clustering with $\hat{k} = 4$



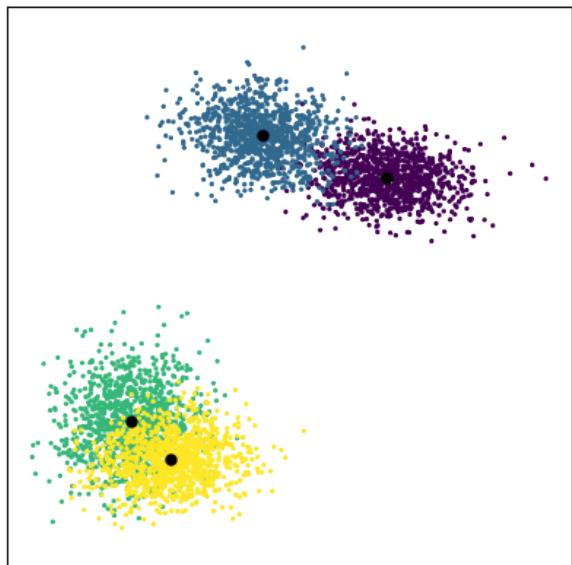
Clustering with $\hat{k} = 5$



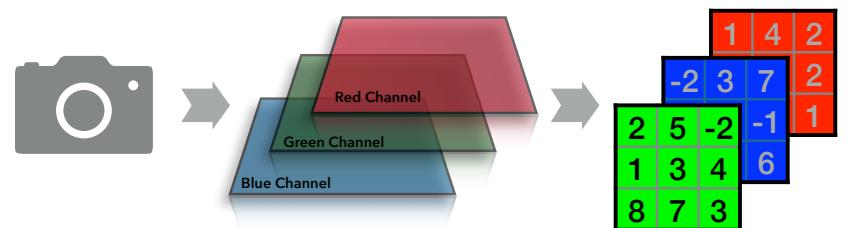
25 / 32

26 / 32

Means of Data Generator



Color Segmentation Using Clustering



- Color images are captured using a CCD device
- Resulting image: 3-dimensional $3 \times n_x \times n_y$ tensor M
- Intensity for color-channel c at location (i, j) in image: $M[c, i, j]$
- $M[c, i, j] \in [p]$ where p ranges from 2^8 ("sdef") to 2^{24} ("true")
- Color palette is $p^3 \gg n_x \times n_y \Rightarrow$ compress image?

27 / 32

28 / 32

Color Segmentation Using Clustering

- View each pixel as a vector in \mathbb{R}^3
- Flatten image, creating a matrix X of $n_x \times n_y$ rows
- Run k-Means on X where $k \ll p^3$
- Resulting k centroids are also colors
- Replace each original pixel with its nearest centroid
- Requires $3pk$ bits to store the exact values of centroids
- Each pixel is replaced with index ($\log(k)$ bits) of its centroid
- Number of bits:

Compressed	$3pk + n_x n_y \log(k)$
Original	$n_x n_y p^3$

- Moreover, images are naturally segmented by color

▶ Code (in → port 9999)

```
In [1]: # General imports
import urllib
import matplotlib.pyplot as plt
from matplotlib.image import imread
from sklearn.cluster import KMeans

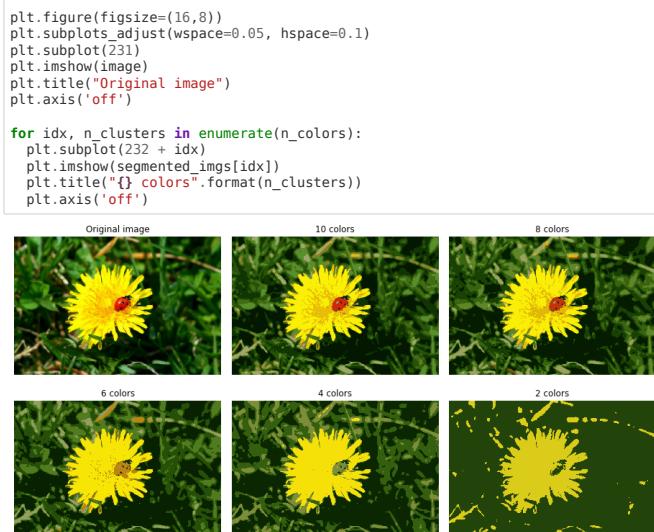
In [2]: # Data Loading
droot = "https://raw.githubusercontent.com/ageron/handson-ml2/master/"
filename = "ladybug.png"
url = droot + "images/unsupervised_learning/" + filename
with urllib.request.urlopen(url) as f:
    image = imread(f)
    X = image.reshape(-1, 3)

In [3]: # Clustering using sklearn.KMeans and then using centroids for segmentation
segmented_imgs = []
n_colors = (10, 8, 6, 4, 2)
for n_clusters in n_colors:
    kmeans = KMeans(n_clusters=n_clusters, random_state=17).fit(X)
    segimg = kmeans.cluster_centers_[kmeans.labels_]
    segmented_imgs.append(segimg.reshape(image.shape))
```

29 / 32

30 / 32

In [4]: # Plot images replacing each pixel with each nearest centroid



Finished $\not\Rightarrow$ Complete

*When you learn from the wrong person you are finished.
When you learn from the right person you are complete.*

[circa 1995]

- Decision trees
- Temporal Modeling
(Markov Models & Random Fields, Recurrent Networks)
- Ensemble Methods & Boosting
- Dimensionality Reduction
- Planning & Markov Decision Processes

31 / 32

32 / 32