

CLUSTERING

- Unsupervised learning problem

- Unlabelled data  $\mathcal{D} = \{\underline{x}_i\}_{i=1}^N$

- Goal: To divide  $\mathcal{D}$  into clusters (e.g. to create labels)

K-means clustering

- Use Euclidean distance as metric

- Minimize

$$J_K(\{\underline{\mu}_l\}_{l=1}^K, \{\underline{z}_i\}_{i=1}^N) = \sum_{i=1}^N \|\underline{x}_i - \underline{\mu}_{\underline{z}_i}\|^2$$

- $\underline{\mu}_l$  is "mean" of cluster  $l$ ,  $l = 1, \dots, K$

- $\underline{z}_i$  is index of cluster given to datapoint  $\underline{x}_i$   
 $\underline{z}_i$  takes values in  $\{1, 2, \dots, K\}$

- Exact minimization of  $J_K$  difficult

- Iterative algorithm to approximate minimization

- Initialize  $\{\underline{\mu}_l\}_{l=1}^K$  to  $K$  random points

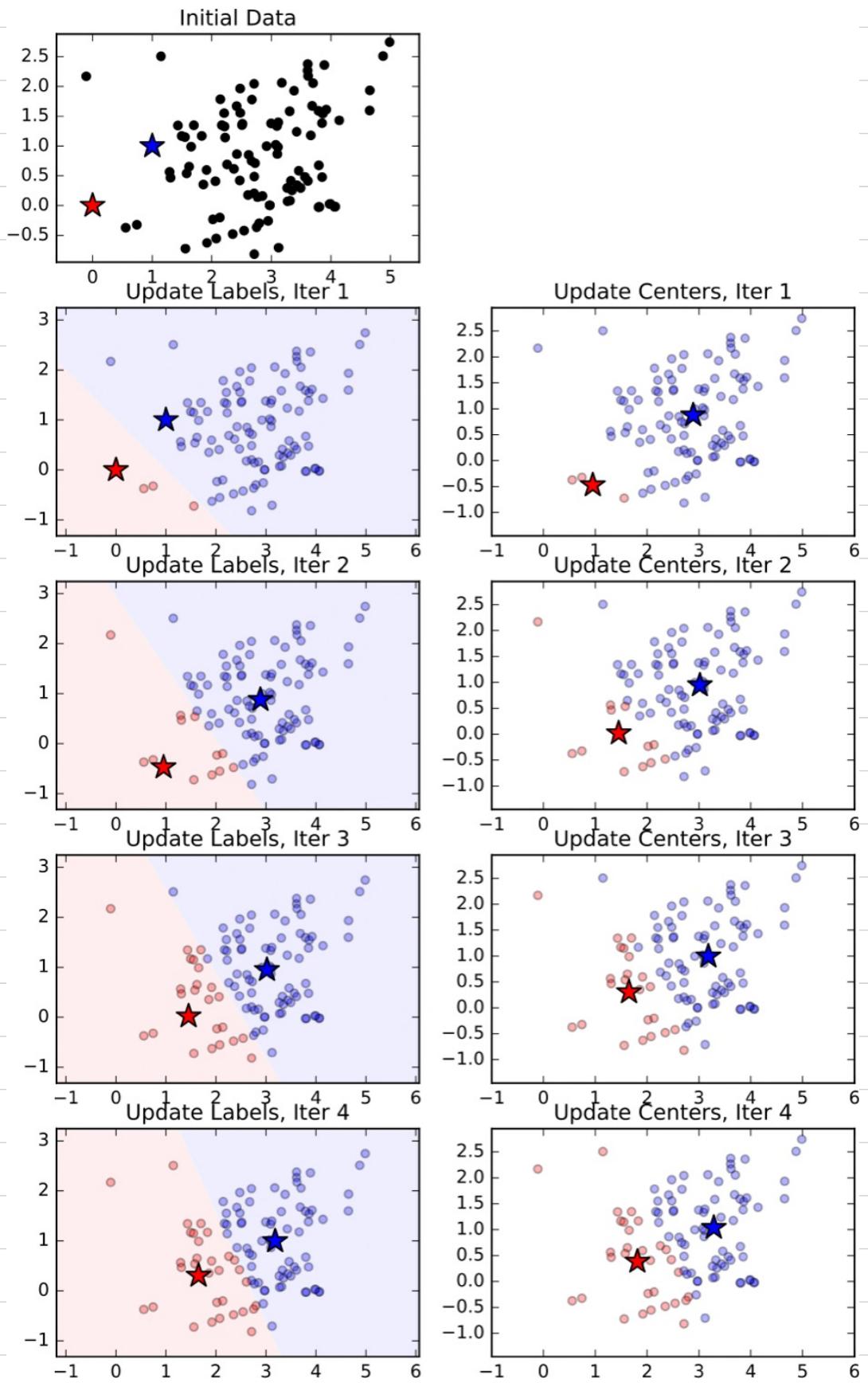
- Assign  $\underline{x}_i$  to cluster index with closest mean

i.e.  $\underline{z}_i = \arg \min_{l=1, \dots, K} \|\underline{x}_i - \underline{\mu}_l\|^2$

- After all data assigned, recompute  $\{\underline{\mu}_l\}_{l=1}^K$

$$\underline{\mu}_l = \frac{1}{N_l} \sum_{i: \underline{z}_i=l} \underline{x}_i, \quad l=1, 2, \dots, K$$

- Repeat until convergence ( $\underline{\mu}_l$ 's don't change)



## K-Medoids Algorithm

- Replace Euclidean distance  $\|\underline{x}_i - \underline{\mu}_e\|$  with more general dissimilarity measure  $d(\underline{x}_i, \underline{\mu}_e)$
- Replace "mean" with "centroid" (or "mediod")

## How to pick K?

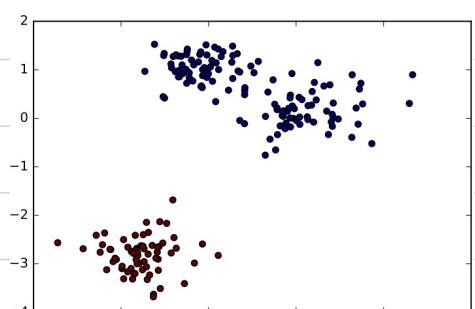
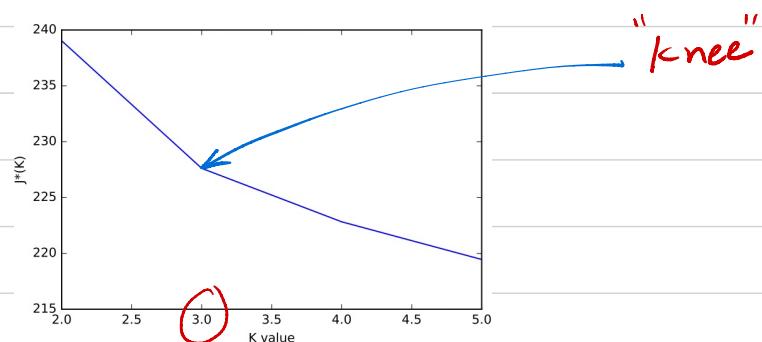
- Sometimes  $K$  is known beforehand, e.g., for handwritten digit recognition  $K = 10$
- More generally, need to estimate  $K$  from data

Use K-means objective  $J^*(K)$

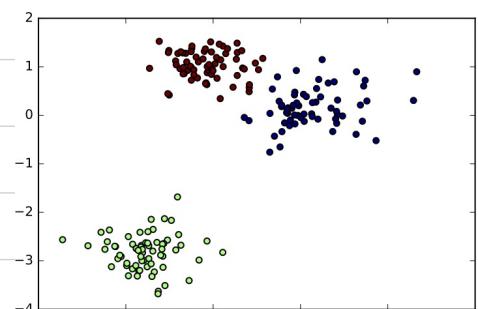
$$J^*(K) = \min_{\{\underline{z}_i\}_{i=1}^N, \{\underline{\mu}_e\}_{e=1}^K} \sum_{i=1}^N \|\underline{x}_i - \underline{\mu}_{z_i}\|^2$$

- Estimate  $J^*(K)$  by running K-means algorithm with many starting points and finding the smallest resultant metric

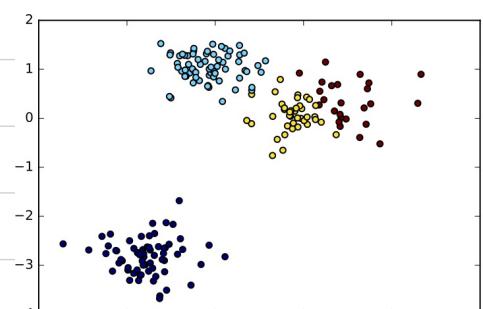
- Plot  $J^*(K)$  v.  $K$  and look for "knee"



$K = 2$



$K = 3$



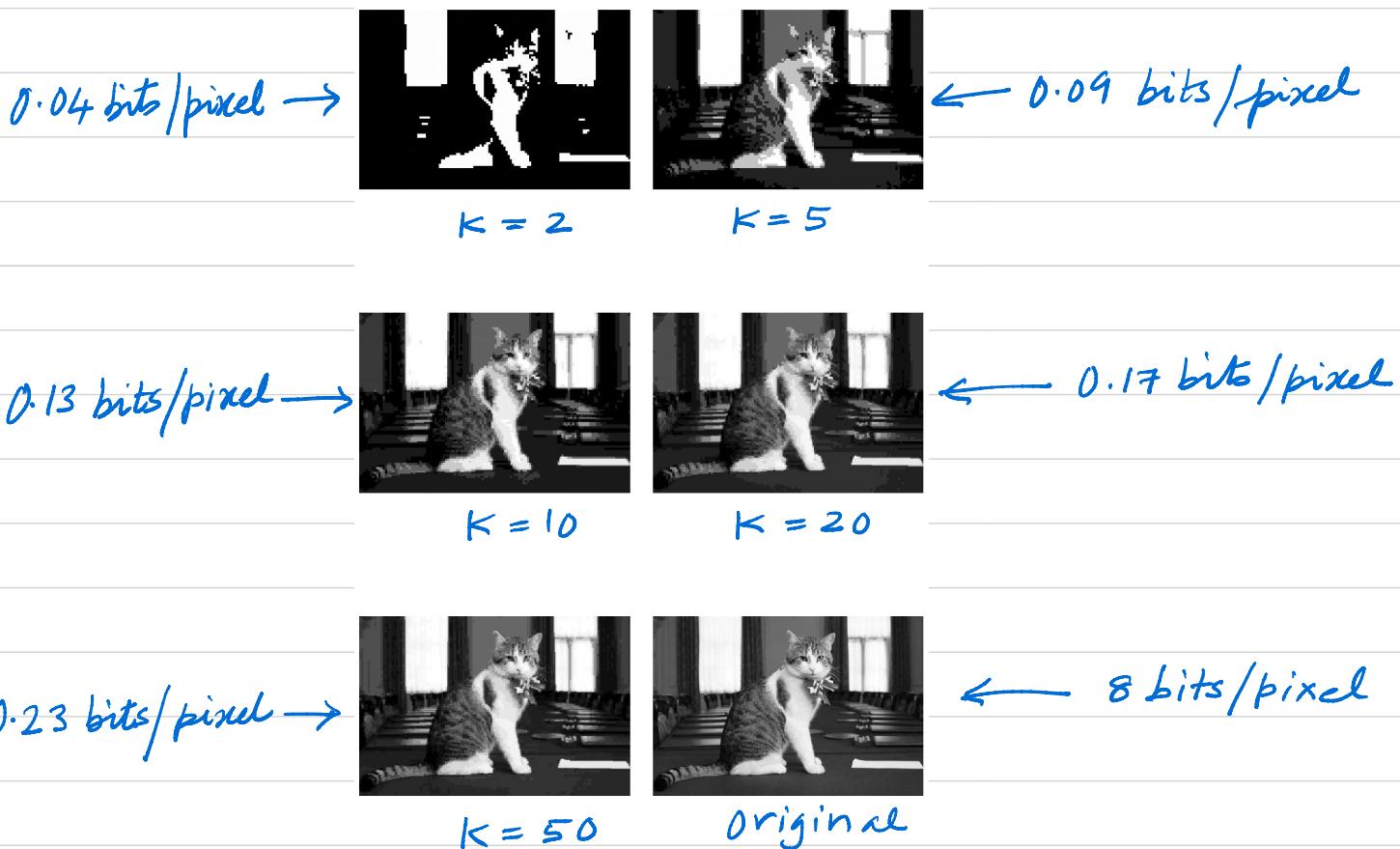
$K = 4$

## Applications of Clustering

Obvious application: creating labelled data from unlabelled data for classification

## Vector Quantization (VQ)

- Lossy compression of images
- chop image up into equal size blocks  
e.g.  $5 \text{ pixels} \times 5 \text{ pixels} = 25\text{-dim vector } \underline{x}_i$
- Do K-means clustering of  $\underline{x}_i$ 's for some K
- Encode image by storing the K means  $\{\underline{\mu}_k\}_{k=1}^K$  and cluster indices  $\{z_i\}_{i=1}^N$  for the  $\underline{x}_i$ 's
- To decompress image use the mean  $\underline{\mu}_{z_i}$  in place of  $\underline{x}_i$



## Image Segmentation

- Similar to VQ

$$\underline{x}_i = [R \ G \ B]$$

Values for each pixel  $i$

- Cluster  $\underline{x}_i$ 's using k-means, and replace pixel with mean  $\mu_{zi}$



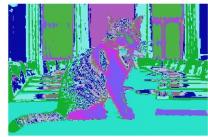
$K = 2$



$K = 3$



$K = 5$



$K = 10$



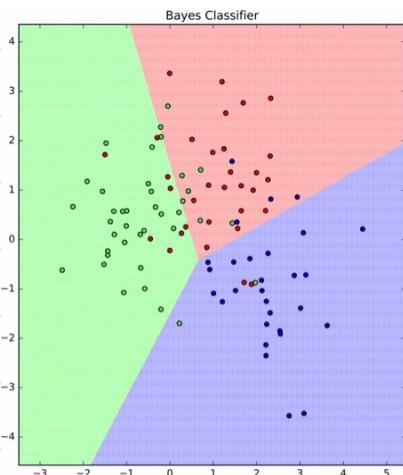
$K = 20$



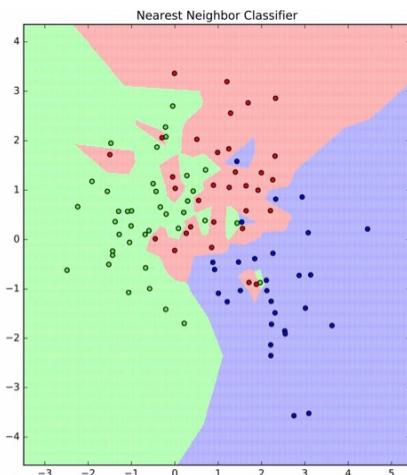
original

## K-Means for NN Classification

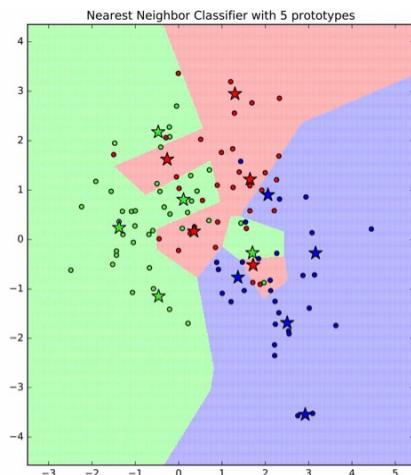
- NN requires computation of  $N$  distances for each new  $\underline{x}$  that we want to classify using  $T = \{\underline{x}_i, y_i\}_{i=1}^N$
- For each label  $l$ , use K-means to cluster the training data for label  $l$  ( $K \ll N$ )
- Replace training data for each class with the k means  $\mu_1, \dots, \mu_K$  for that class.
- Now we need to only compute  $KM$  distances for each new  $\underline{x}$ , which may be  $\ll N$  if  $N$  is large



Bayes Classifier



Nearest Neighbor Classifier



Nearest Neighbor Classifier with 5 prototypes

$N=100, M=3$

$K=5$