

CPSC 340: Machine Learning and Data Mining

K-Means Clustering

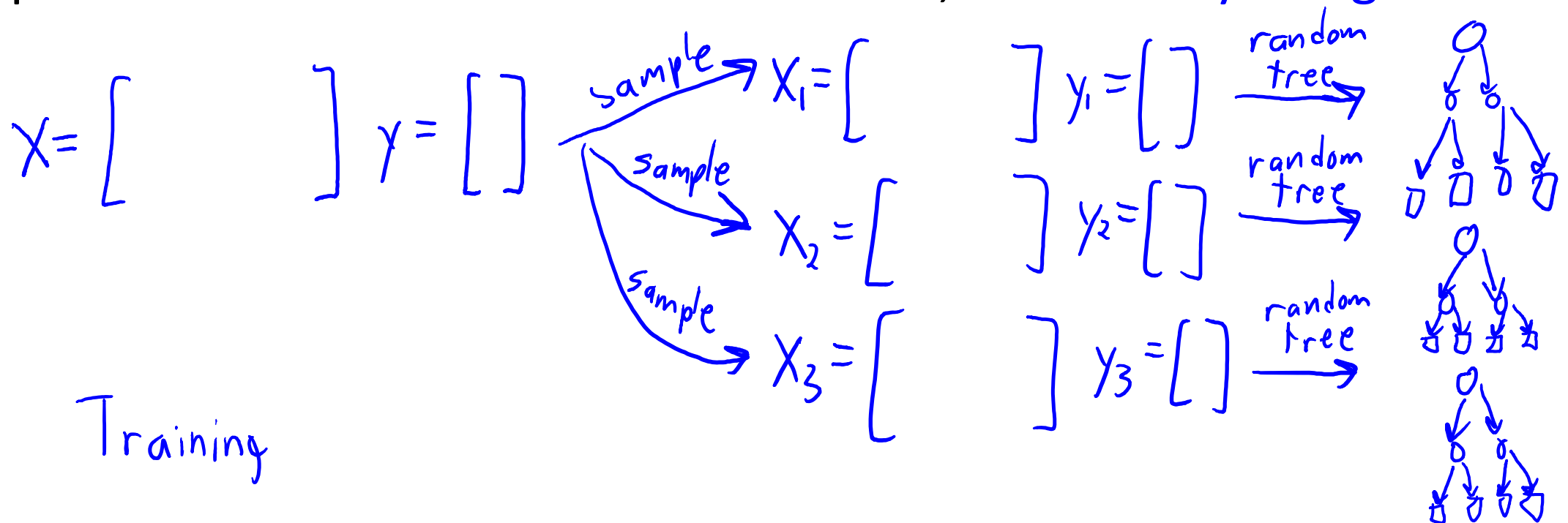
Fall 2016

Admin

- **Assignment 1 is due now!**
 - 1 late day to hand it in before Monday's class.
 - 2 late days to hand it in before Wednesday's class.
 - 3 late days to hand it in before Friday of next week's class.
 - 0 after that.
- Assignment 2 coming next week.

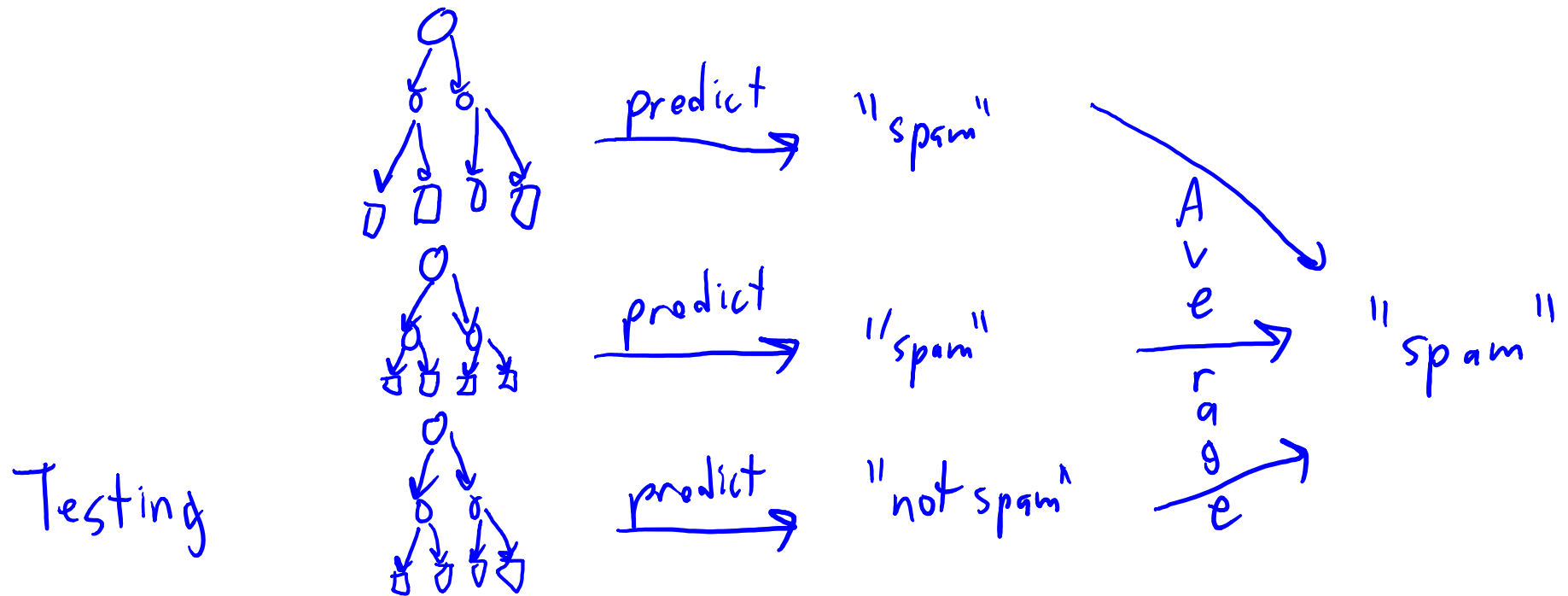
Random Forests

- Random forests are one of the best 'out of the box' classifiers.
- Fit deep decision trees to random **bootstrap samples** of data, base splits on **random subsets** of the features, and **classify using mode**.



Random Forests

- Random forests are one of the best 'out of the box' classifiers.
- Fit deep decision trees to random **bootstrap samples** of data, base splits on **random subsets** of the features, and **classify using mode**.



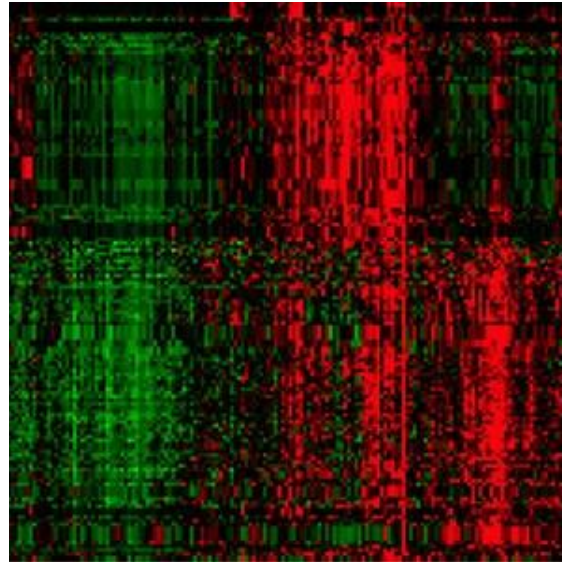
End of Part 1: Key Concepts

- Fundamental ideas:
 - Training vs. test error.
 - Golden rule of ML.
 - Fundamental trade-off.
 - Validation sets and cross-validation.
 - Parametric vs. non-parametric.
 - No free lunch theorem.
 - Ensemble methods.
- Methods that we focused on:
 - Decision trees (greedy recursive splitting using decision stumps).
 - Naïve Bayes (generative classifier based on conditional independence).
 - K-nearest neighbours (non-parametric classifier with universal consistency).
 - Random forests (averaging plus randomization to reduce overfitting).

Application: Classifying Cancer Types

- “I collected gene expression data for 1000 different types of cancer cells, can you tell me the different classes of cancer?”

$X =$



- We are not given the class labels y , but want **meaningful labels**.
- An example of **unsupervised learning**.

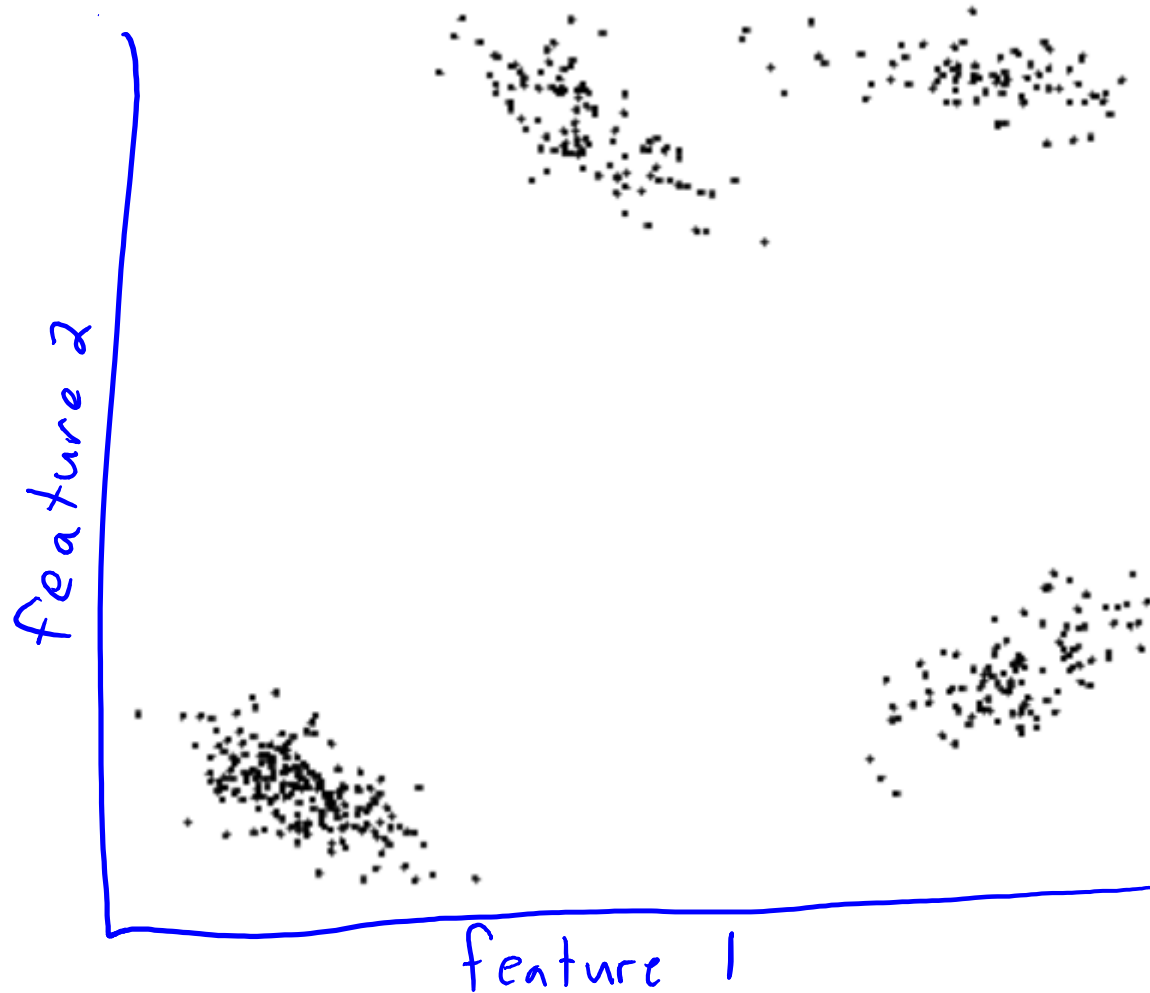
Unsupervised Learning

- Supervised learning:
 - We have features x_i and class labels y_i .
 - Write a program that produces y_i from x_i .
- Unsupervised learning:
 - We **only have x_i values**, but no explicit target labels.
 - You want to do ‘something’ with them.
- Some unsupervised learning tasks:
 - Outlier detection: Is this a ‘normal’ x_i ?
 - Data visualization: What does the high-dimensional X look like?
 - Association rules: Which x_{ij} occur together?
 - Latent-factors: What ‘parts’ are the x_i made from?
 - Ranking: Which are the most important x_i ?
 - Clustering: What types of x_i are there?

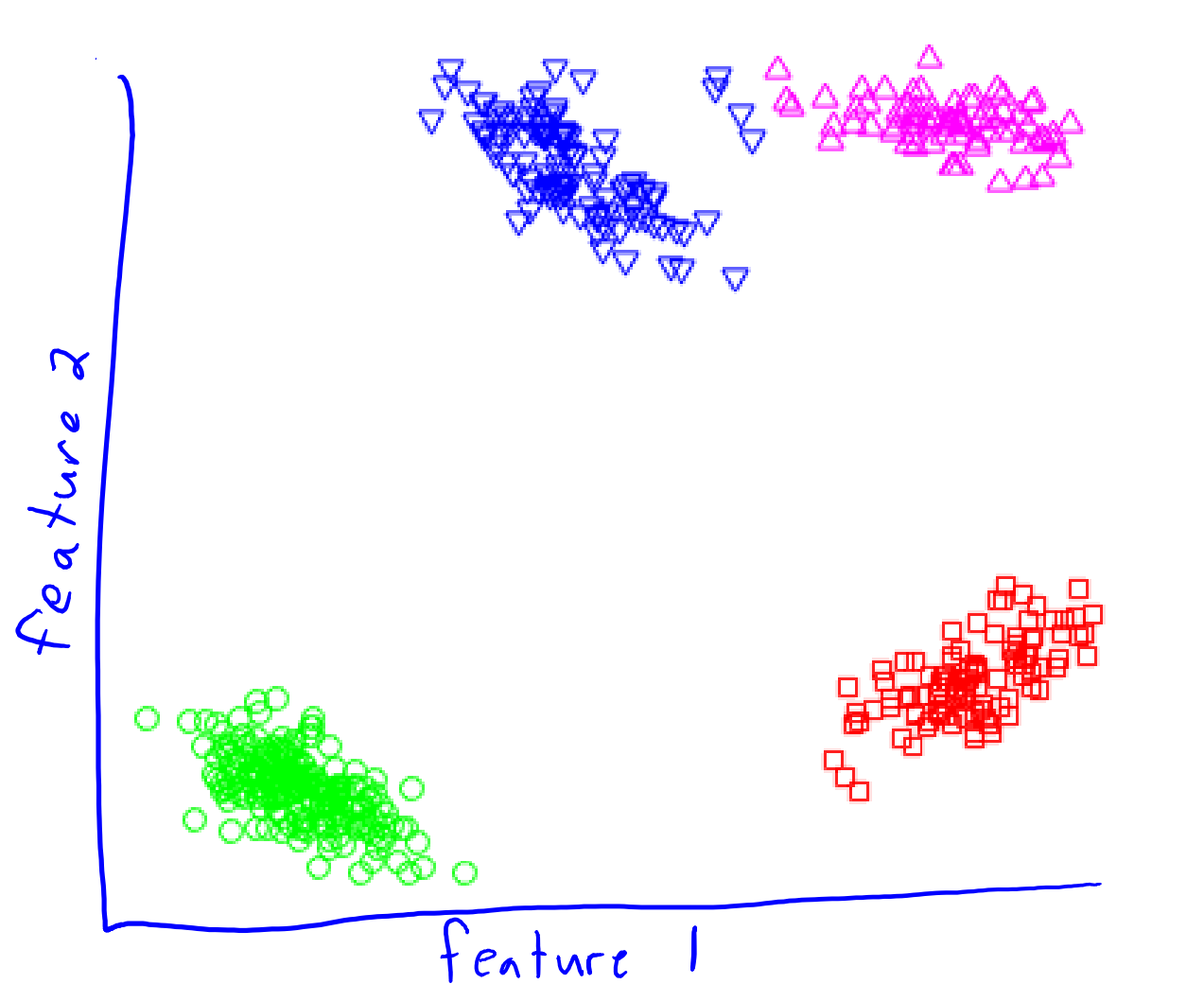
Clustering

- Clustering:
 - Input: set of objects described by features x_i .
 - Output: an assignment of objects to 'groups'.
- Unlike classification, we are not given the 'groups'.
 - Algorithm must discover groups.
- Example of groups we might discover in e-mail spam:
 - 'Lucky winner' group.
 - 'Weight loss' group.
 - 'Nigerian prince' group.

Clustering Example



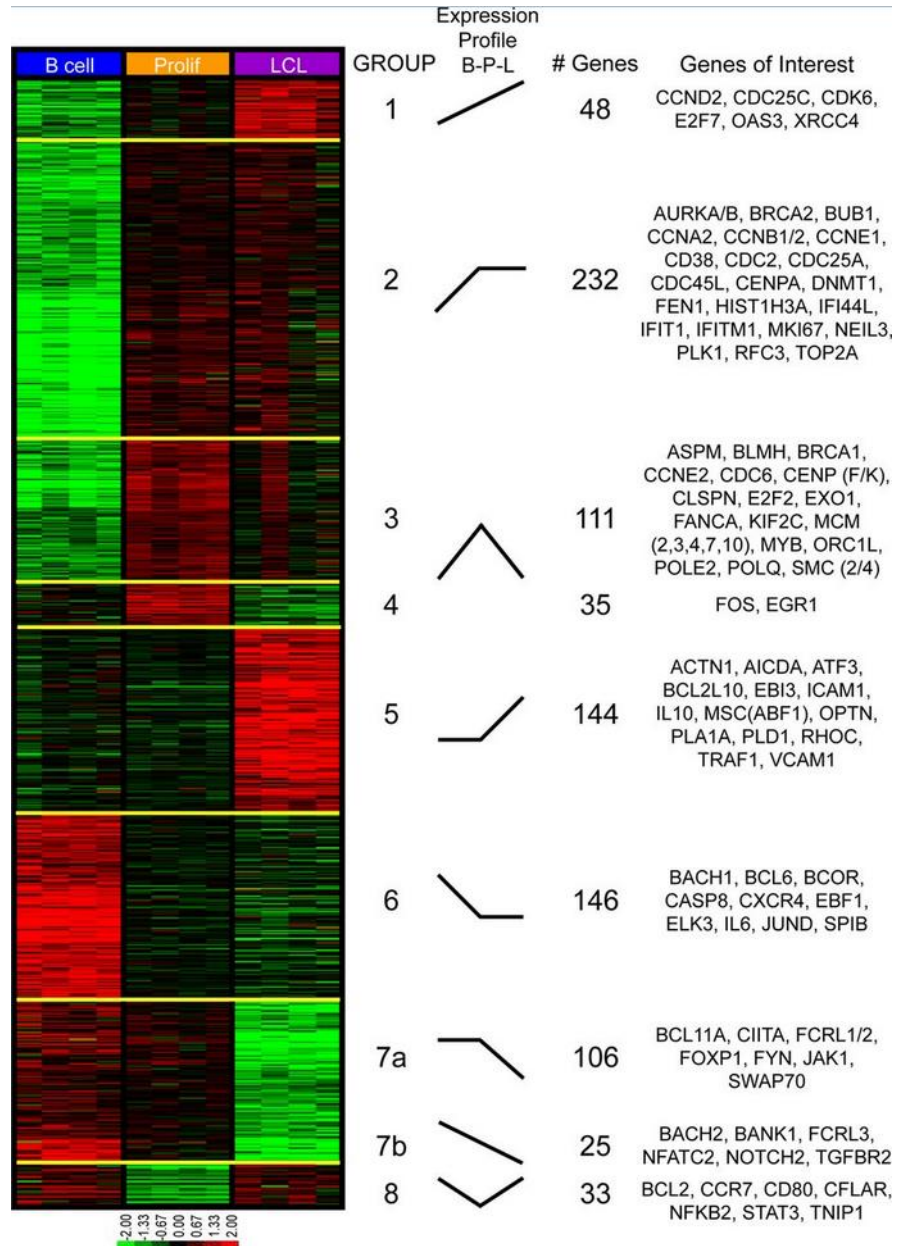
Clustering Example



Data Clustering

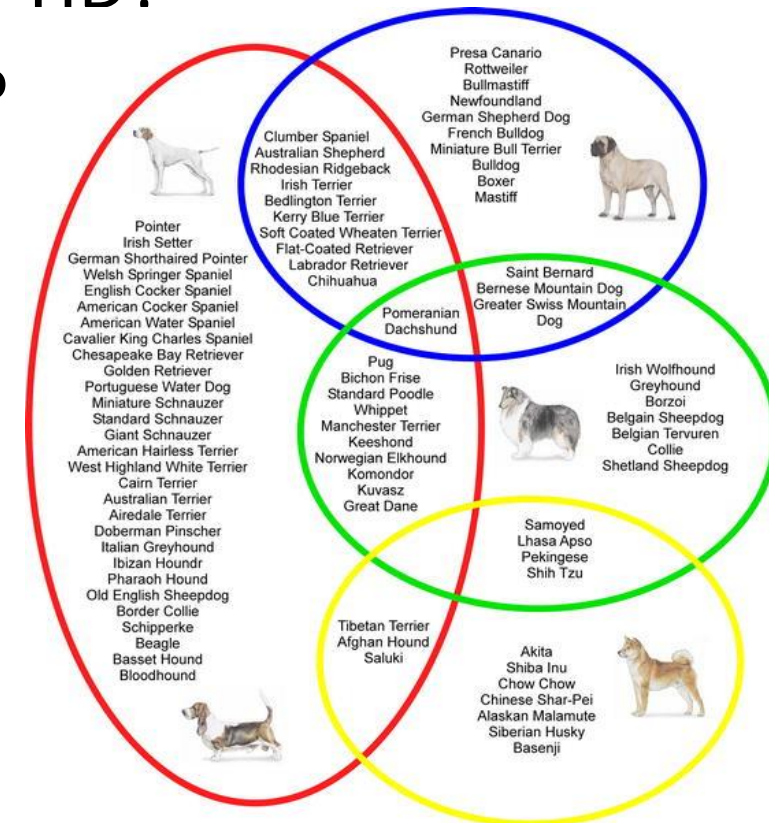
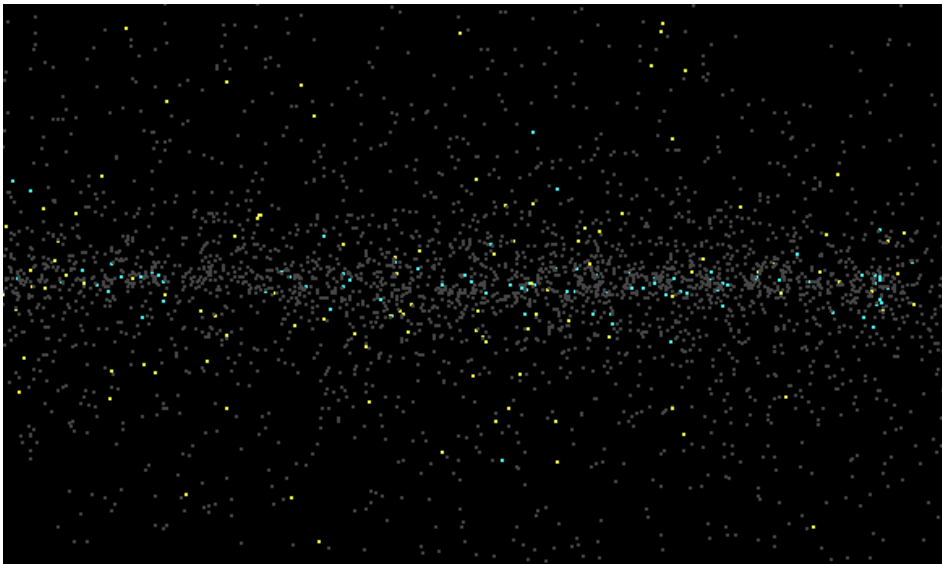
- General goal of clustering algorithms:
 - Objects in the same group should be ‘similar’.
 - Objects in different groups should be ‘different’.
- But the ‘best’ clustering is hard to define:
 - We don’t have a test error.
 - Generally, there is no ‘best’ method in unsupervised learning.
 - Means there are lots of methods: we’ll focus on important/representative ones.
- Why cluster?
 - You could want to know what the groups are.
 - You could want a ‘prototype’ example for each group.
 - You could want to find the group for a new example x .
 - You could want to find objects related to a new example x .

Clustering of Epstein-Barr Virus



Other Clustering Applications

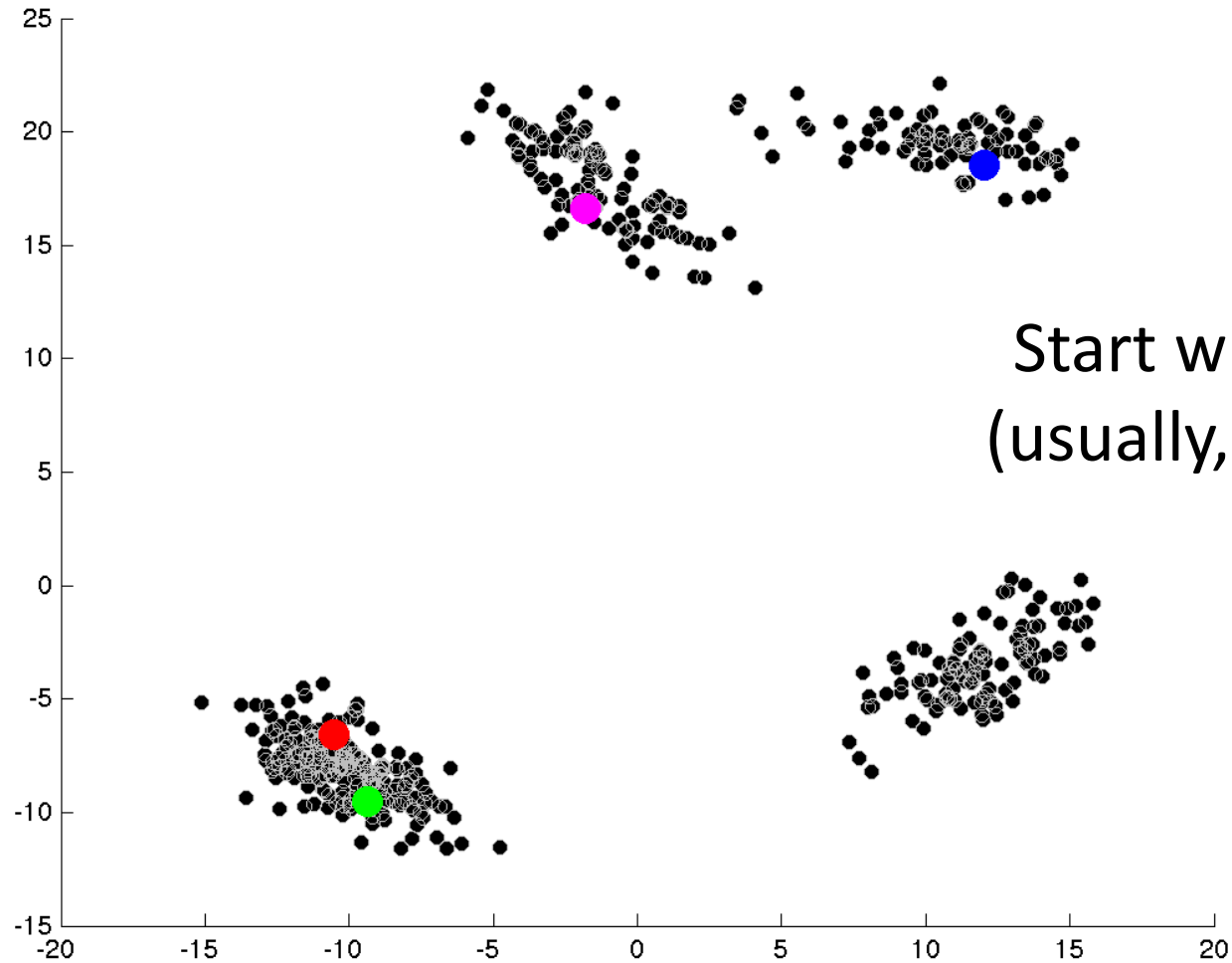
- NASA: what types of stars are there?
- Biology: are there sub-species?
- Documents: what kinds of documents are on my HD?
- Commercial: what kinds of customers do I have?



K-Means

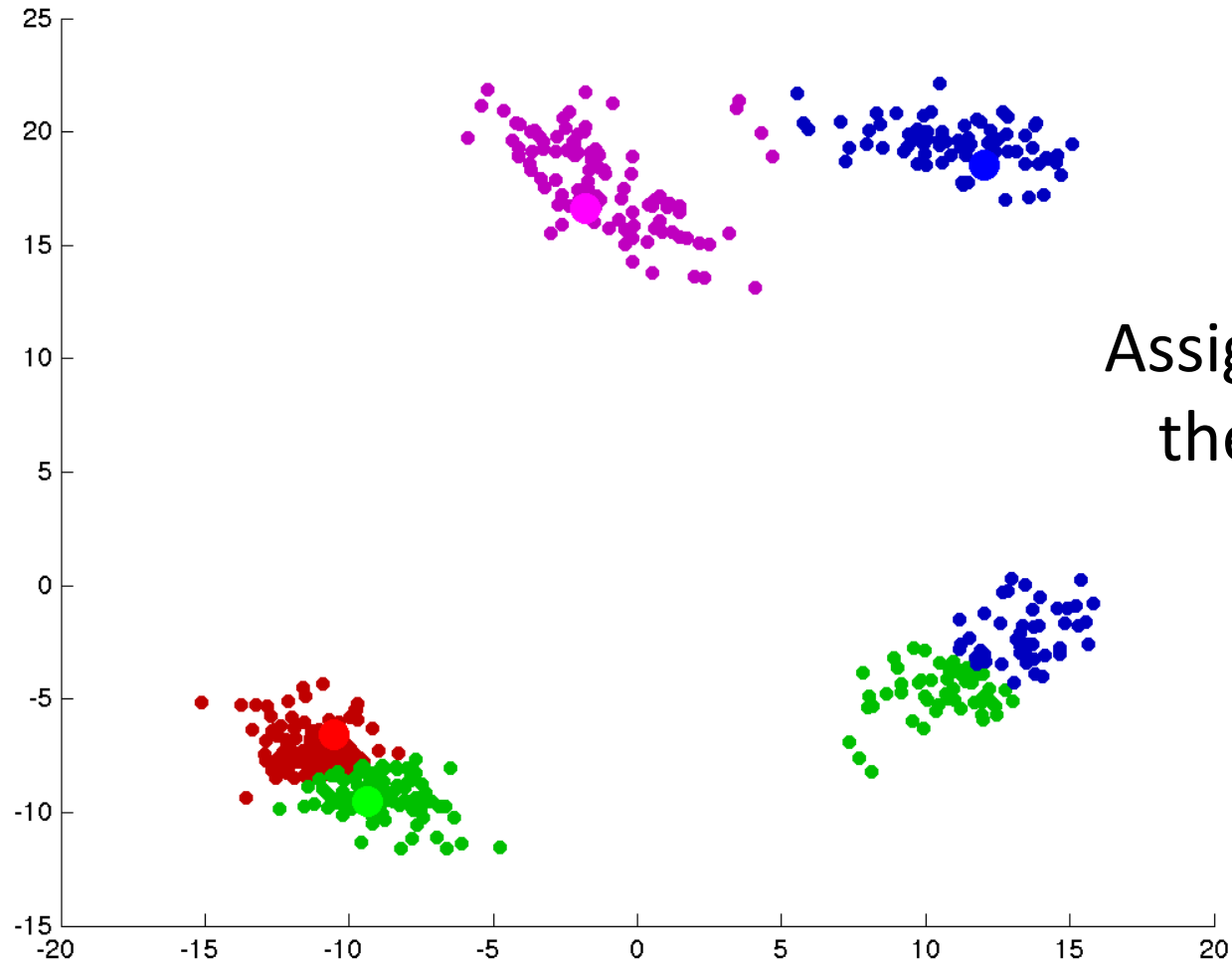
- Most popular clustering method is **k-means**.
- Input:
 - The **number of clusters 'k'**.
 - **Initial guesses of the center ("mean") of each cluster.**
- Algorithm:
 - **Assign each x_i to its closest mean.**
 - **Update the means** based on the assignment.
 - Repeat until convergence.

K-Means Example



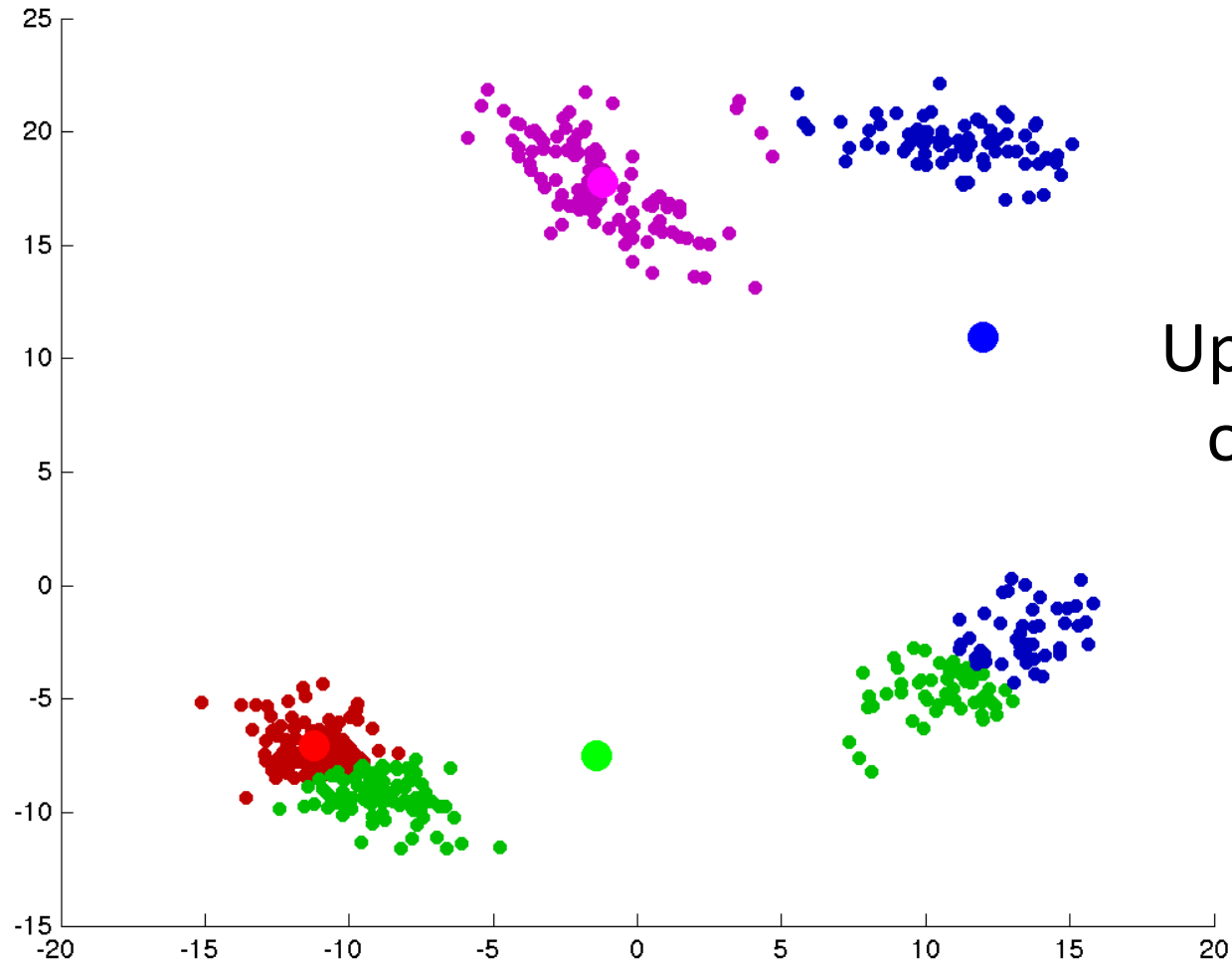
Start with 'k' initial 'means'
(usually, random data points)

K-Means Example



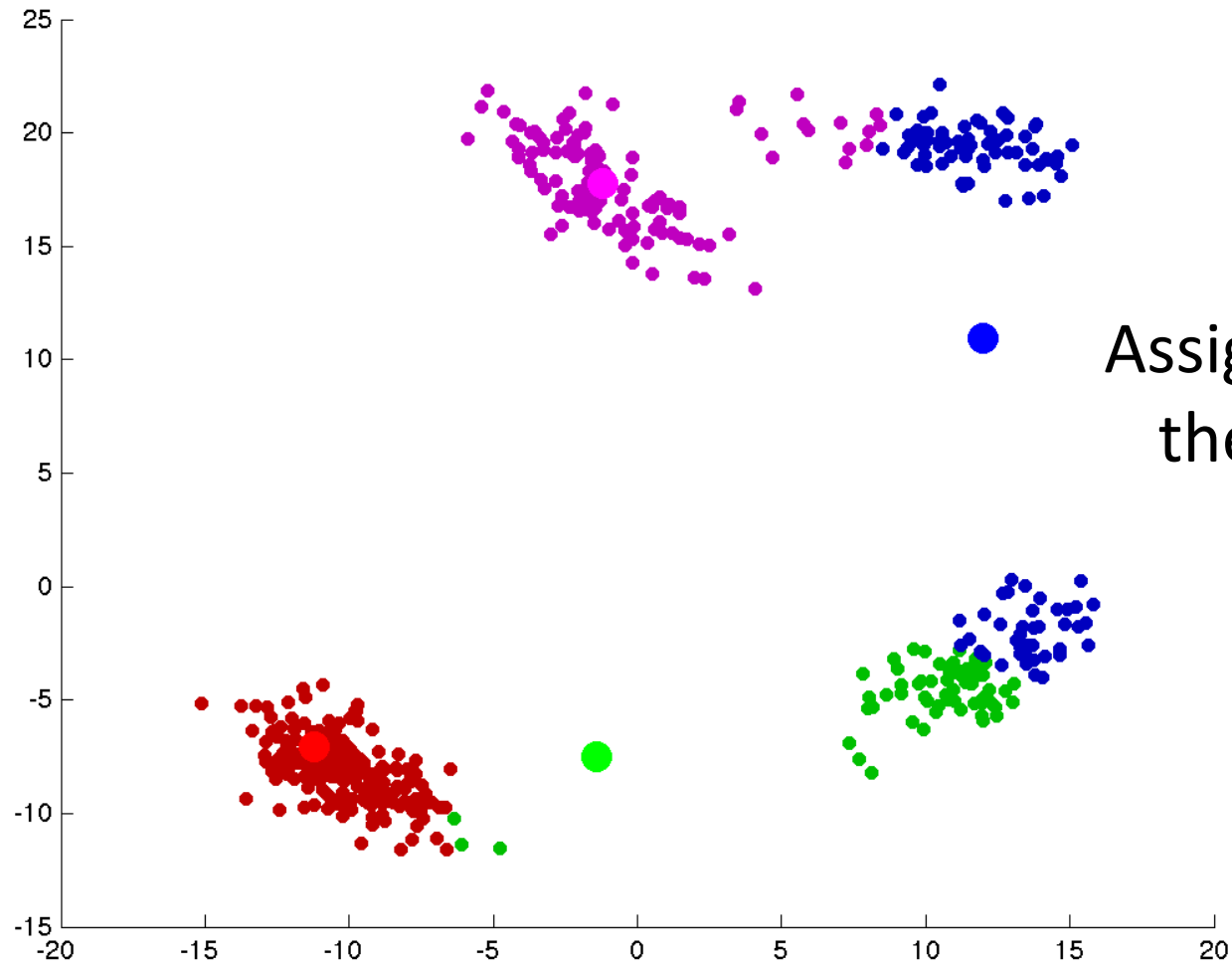
Assign each object to the closest mean.

K-Means Example



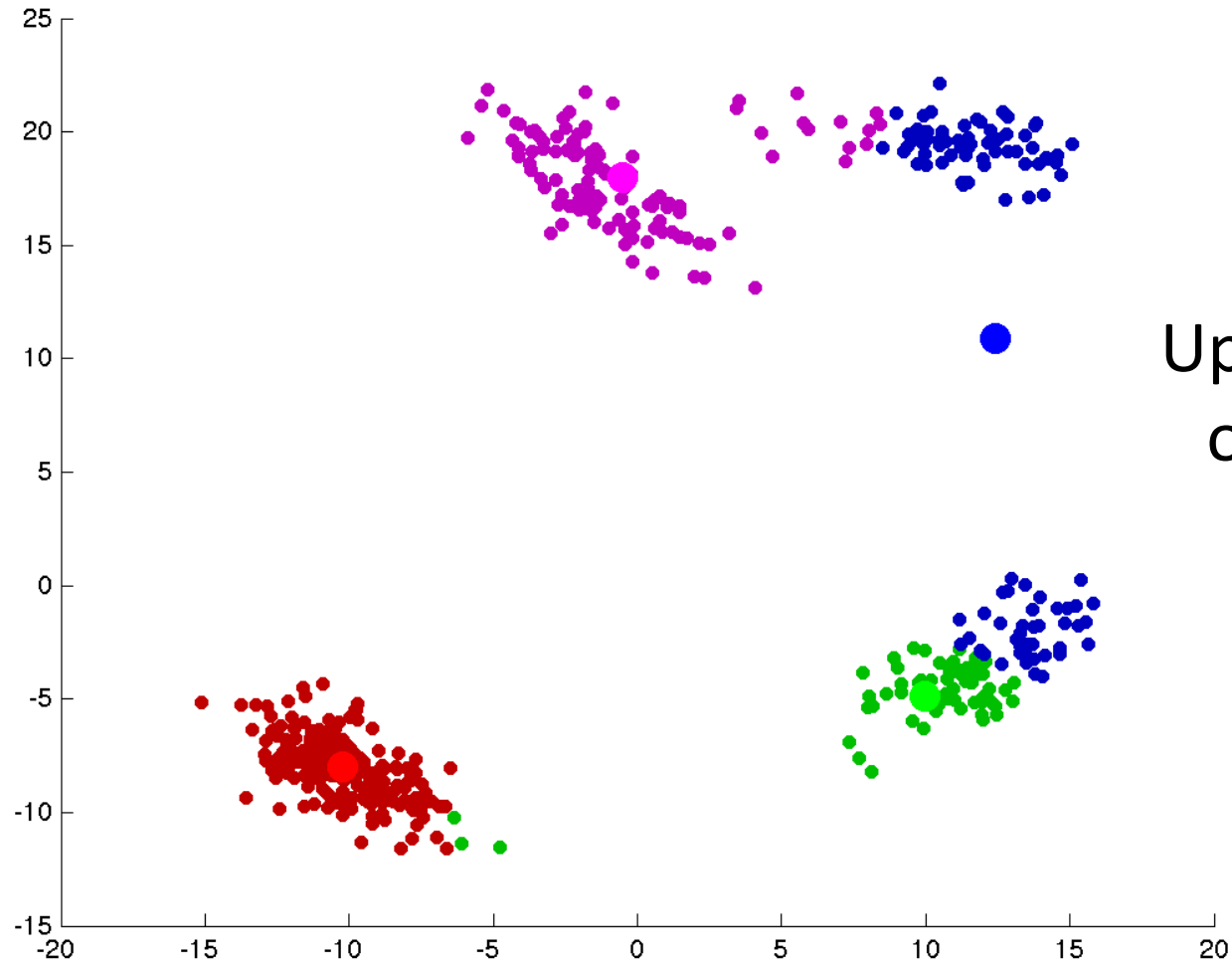
Update the mean
of each group.

K-Means Example



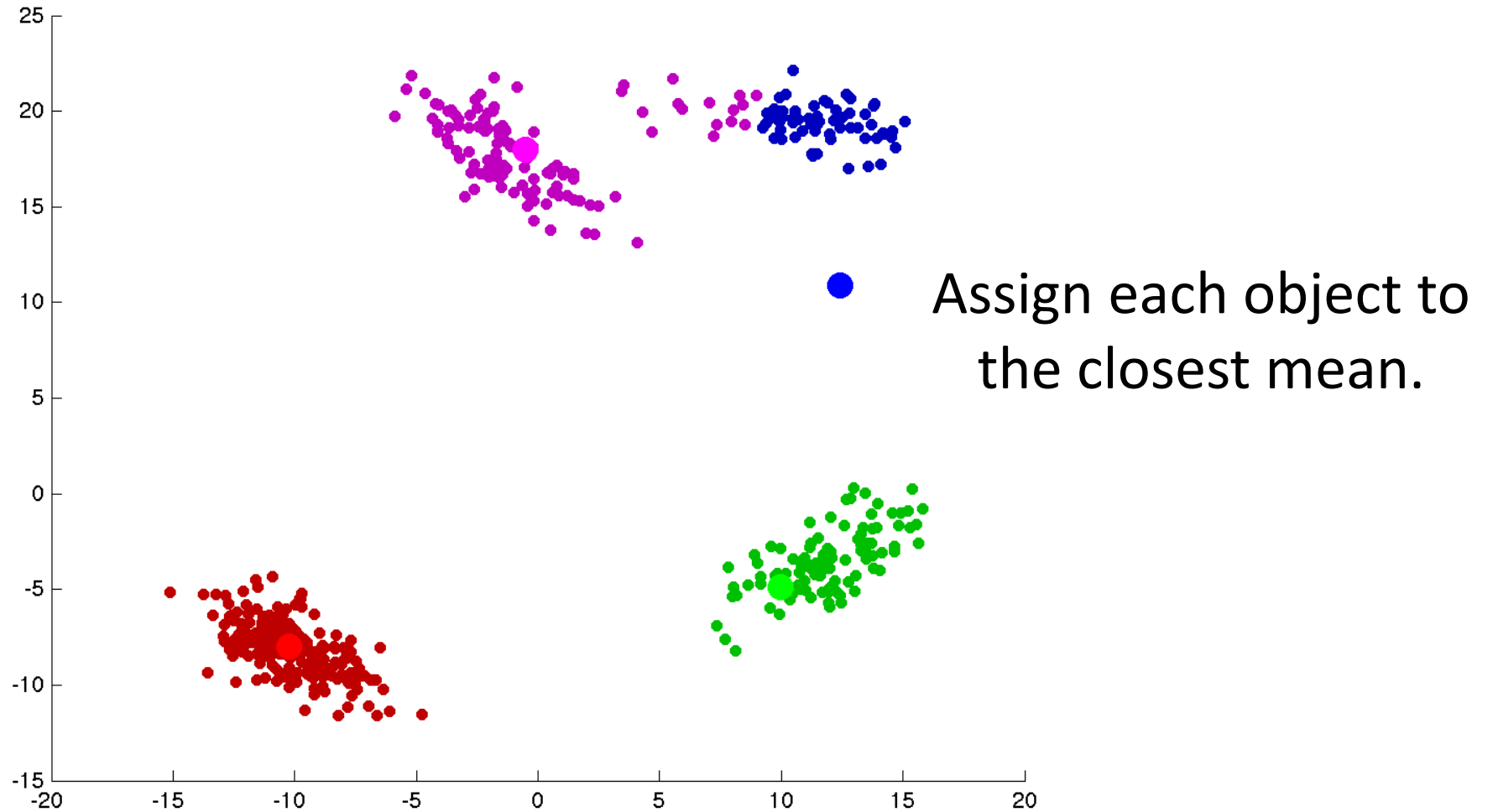
Assign each object to the closest mean.

K-Means Example

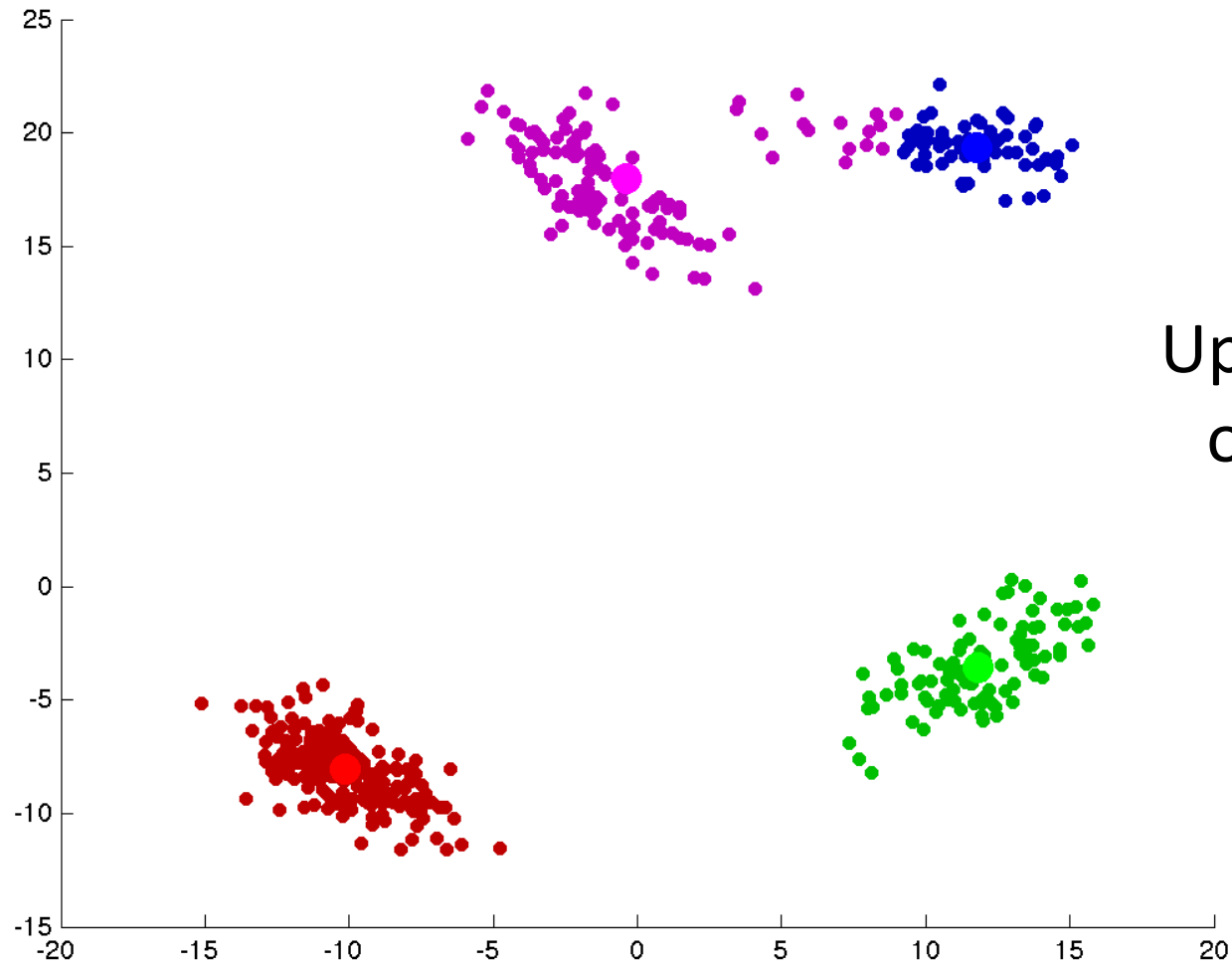


Update the mean
of each group.

K-Means Example

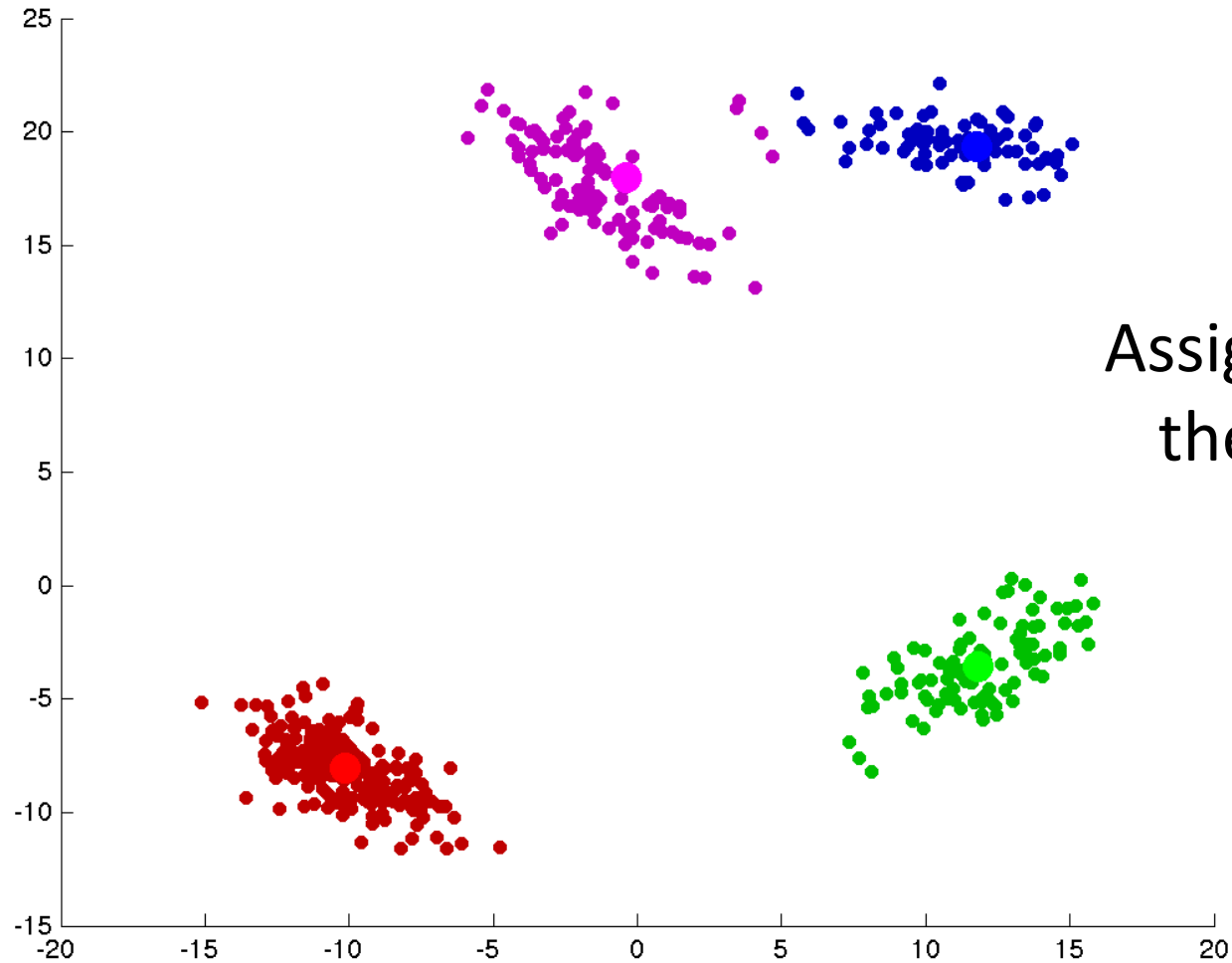


K-Means Example



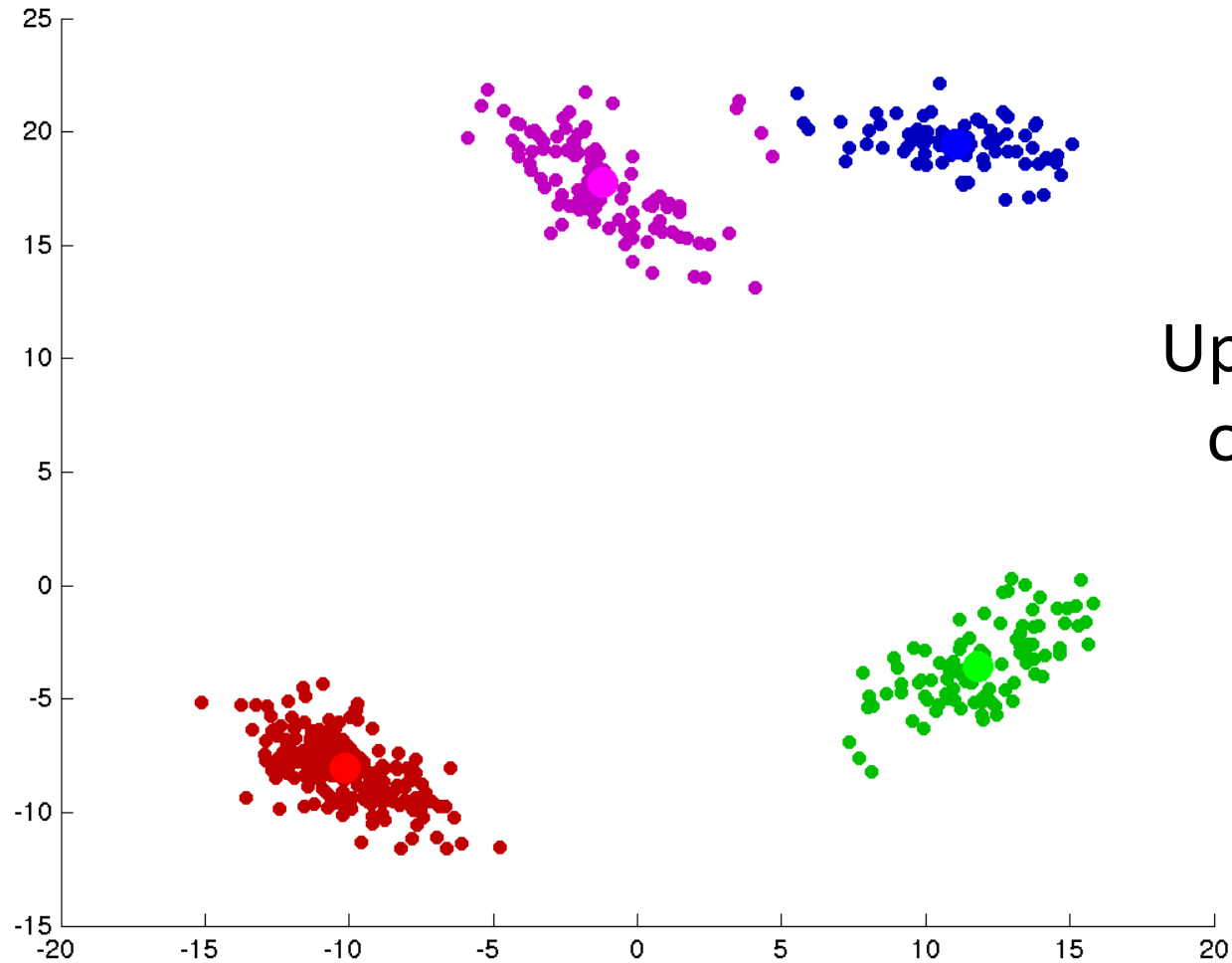
Update the mean
of each group.

K-Means Example



Assign each object to
the closest mean.

K-Means Example



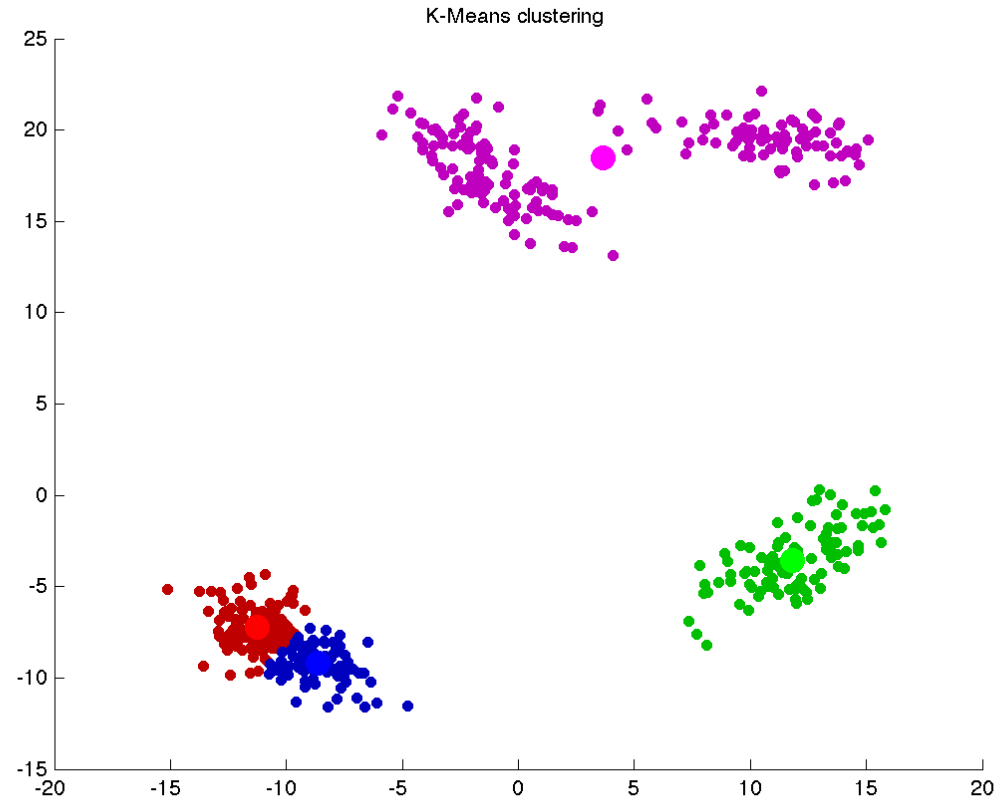
Update the mean
of each group.

Stop if no objects
change groups.

K-Means Issues

- **Guaranteed to converge** when using Euclidean distance.
- **New object are assigned to nearest mean** to cluster them.
- Assumes you **know number of clusters 'k'**.
 - Lots of heuristics to pick 'k', none satisfying:
 - https://en.wikipedia.org/wiki/Determining_the_number_of_clusters_in_a_data_set
- Each object is assigned to **one (and only one) cluster**:
 - No possibility for overlapping clusters or leaving objects unassigned.
- It may converge to **sub-optimal solution...**

K-Means Clustering with Different Initialization



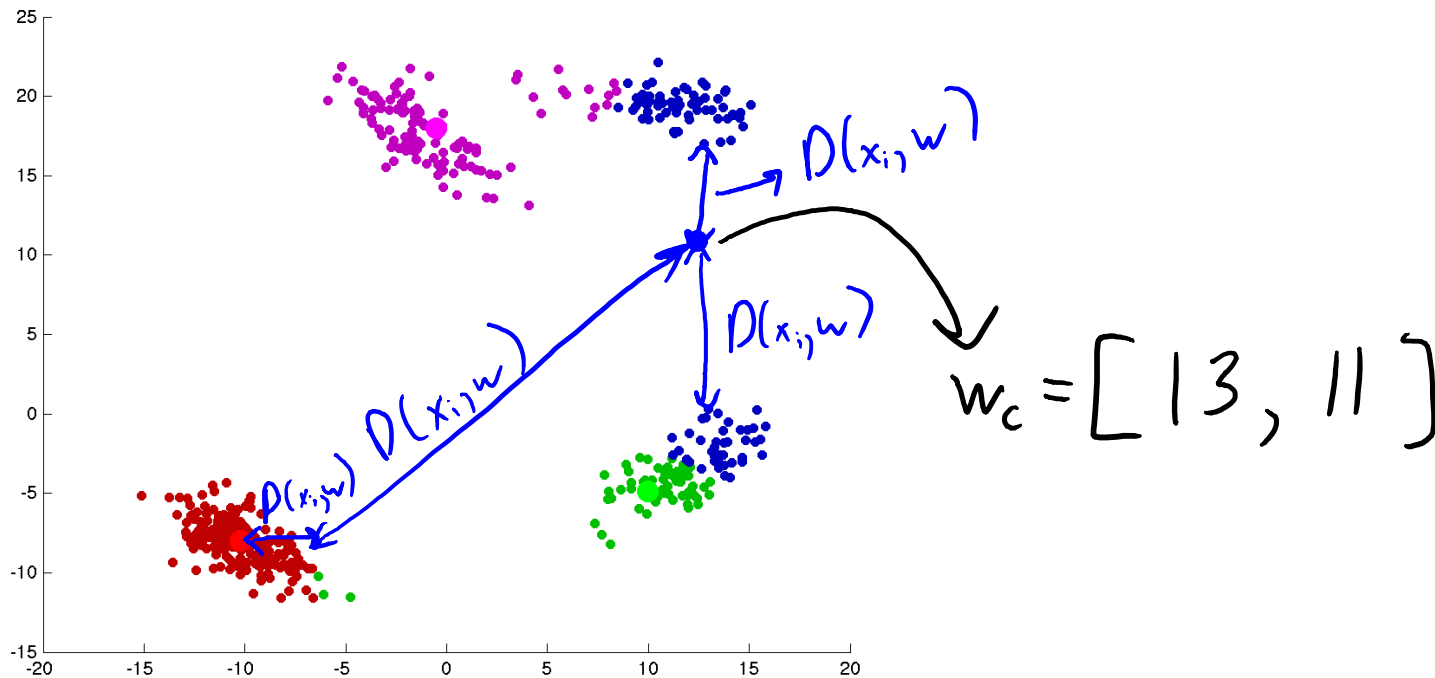
- Classic approach to dealing with sensitivity to initialization:
 - Try several different random starting points, choose the ‘best’.
- We’ll see a more clever approach next time...

Cost of K-means

- Bottleneck is calculating distance from each x_i to each mean w_c :

$$D(x_i, w_c) = \sqrt{\sum_{j=1}^d (x_{ij} - w_{cj})^2}$$

vector
scalar
center of cluster "c"



Cost of K-means

- Bottleneck is **calculating distance** from each x_i to each **mean w_c** :

$$D(x_i, w_c) = \sqrt{\sum_{j=1}^d (x_{ij} - w_{cj})^2}$$

Handwritten annotations:
- A red arrow points from the word "vector" to the w_c term in the equation.
- A red arrow points from the word "center of cluster 'c'" to the w_c term.
- A red arrow points from the word "scalar" to the x_{ij} term.
- A red arrow points from the word "vector" to the x_{ij} term.

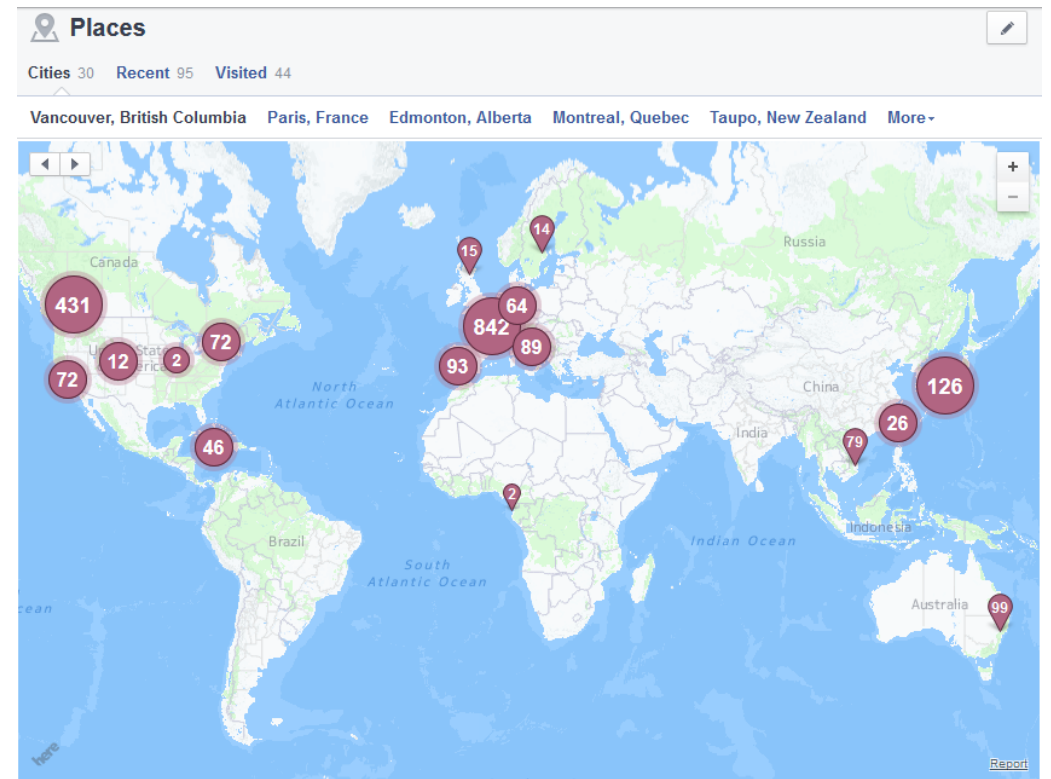
- Each time we do this costs $O(d)$ to go through all features.
- For each of the 'n' objects, we compute the distance to 'k' clusters.
- **Total cost of assigning objects to clusters is $O(ndk)$.**
 - Fast if k is not too large.
- Updating means is cheaper: $O(nd)$.

Handwritten annotation:
- For each cluster 'c', compute $w_c = \frac{1}{n_c} \sum_{i \in C} x_i$
- n_c : Number of objects in cluster 'c'
- $\sum_{i \in C} x_i$: Loop over objects in cluster.
- x_i : Object in cluster.

Vector Quantization

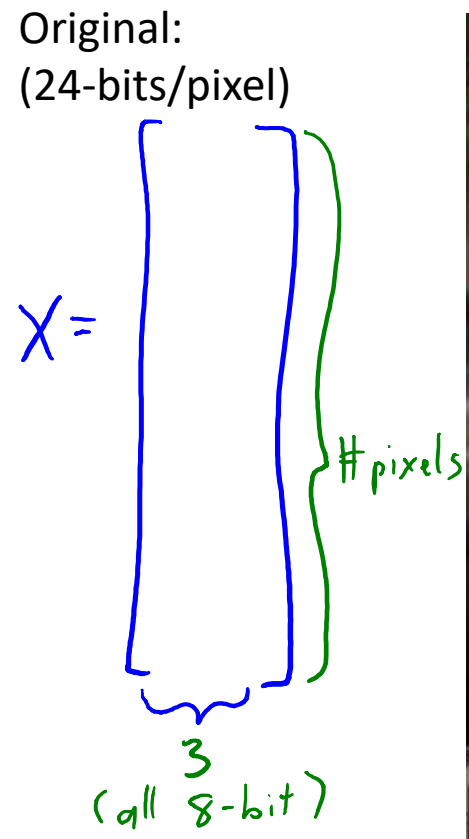
- K-means originally comes from signal processing.
- Designed for **vector quantization**:
 - Replace ‘vectors’ (objects) with a set of ‘prototypes’ (means).

- Example:
 - Facebook places.
 - What sizes of clothing should I make?

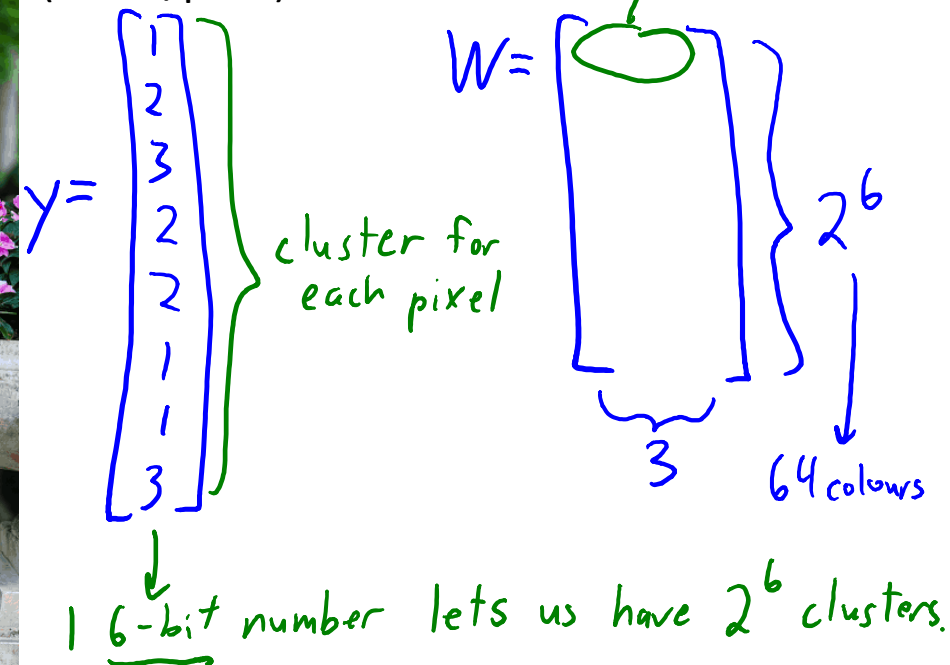


Vector Quantization: Image Colors

- Usual RGB representation of a pixel's color: three 8-bit numbers.
 - For example, [241 13 50] = ■.
 - Can apply k-means to find set of prototype colours.

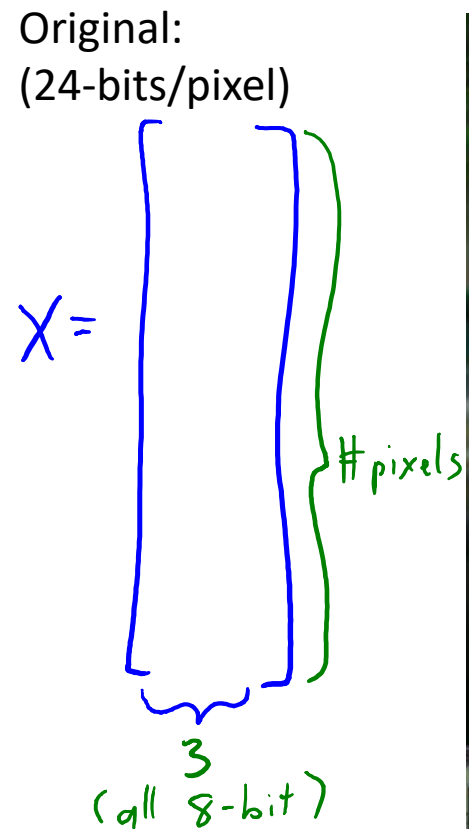


K-Means Quantized:
(6-bits/pixel)

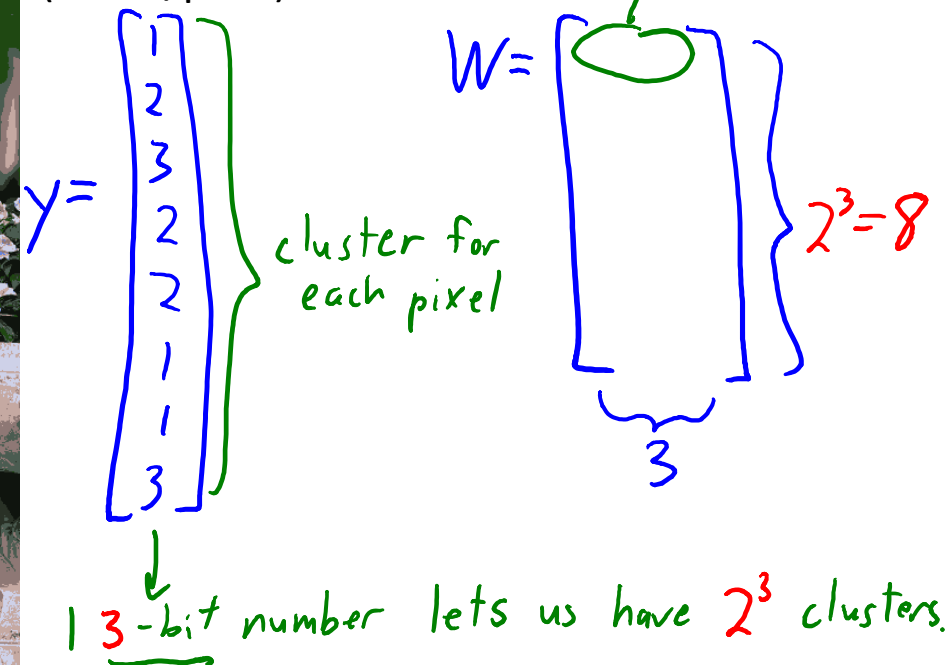


Vector Quantization: Image Colors

- Usual RGB representation of a pixel's color: three 8-bit numbers.
 - For example, [241 13 50] = ■.
 - Can apply k-means to find set of prototype colours.

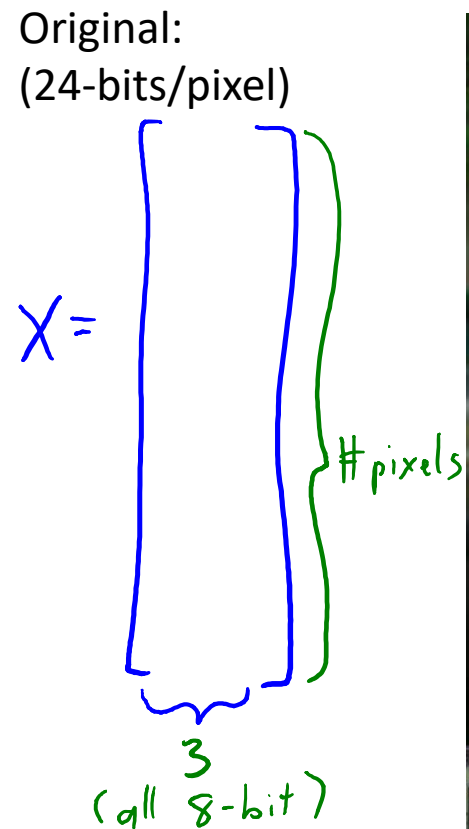


K-Means Quantized:
(3-bits/pixel)

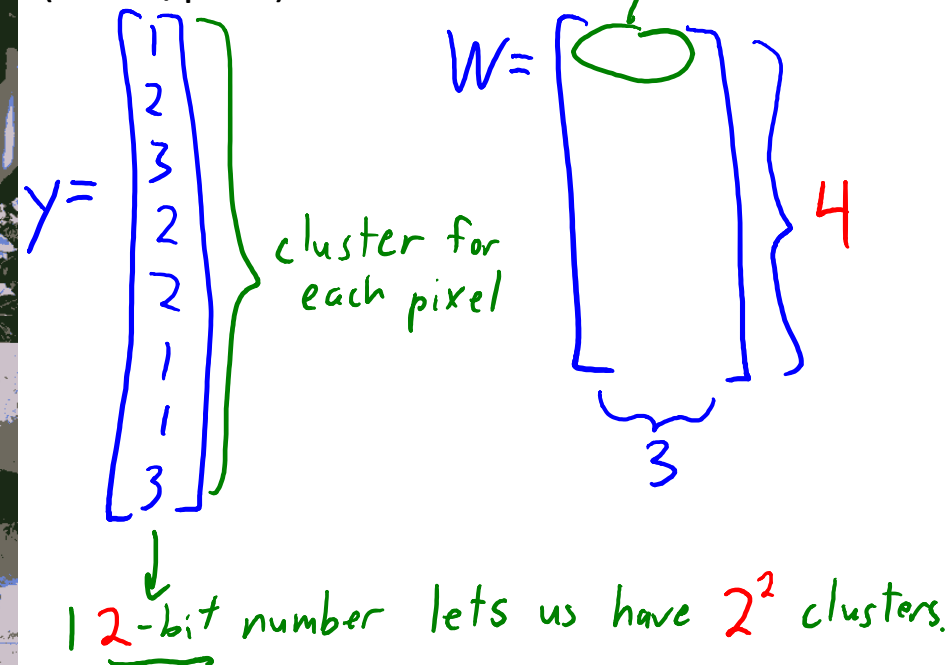


Vector Quantization: Image Colors

- Usual RGB representation of a pixel's color: three 8-bit numbers.
 - For example, [241 13 50] = ■.
 - Can apply k-means to find set of prototype colours.

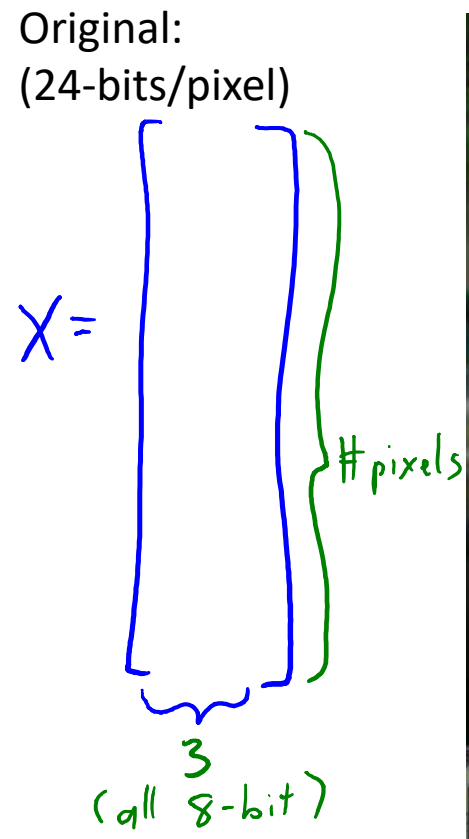


K-Means Quantized:
(2-bits/pixel)

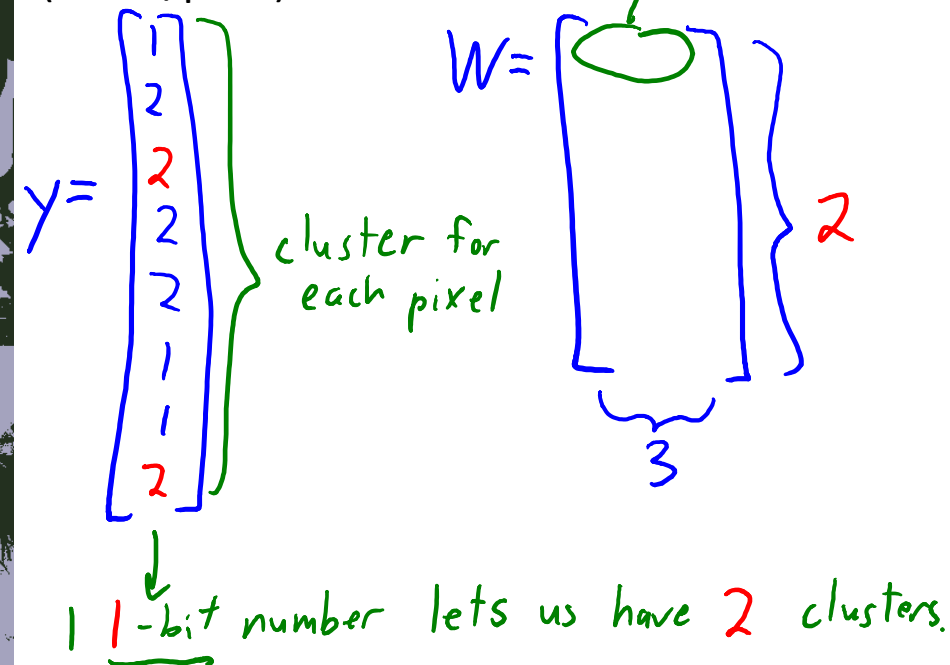


Vector Quantization: Image Colors

- Usual RGB representation of a pixel's color: three 8-bit numbers.
 - For example, [241 13 50] = ■.
 - Can apply k-means to find set of prototype colours.



K-Means Quantized:
(1-bits/pixel)



What is K-Means Doing?

- We can interpret K-Means as trying to minimize an objective:
 - Total sum of **squared distances from object x_i to their centers $w_{c(i)}$** :

$$f(w_1, w_2, \dots, w_k, c(1), c(2), \dots, c(n)) = \sum_{i=1}^n \sum_{j=1}^d (x_{ij} - w_{c(i)j})^2$$

- The k-means steps:

- Optimally update cluster assignments $c(i)$.
- Optimally update means w_c .

- Convergence follows because:

- Each step does not increase the objective.
- There are a finite number of assignments to k clusters.

Cluster of example 'i'

K-Medians Clustering

- With other distances, k-means may not converge.
- However, changing objective function gives convergent algorithms.

- E.g., we can use the L1-norm:
$$\sum_{i=1}^n \sum_{j=1}^d |x_{ij} - w_{c(i)j}|$$

- A 'k-medians' algorithm based on the L1-norm:
 - Cluster assignment based on the L1-norm (nearest median).
 - Update 'medians' as median value (dimension-wise) of each cluster.
- This approach is more robust to outliers.

↑ k-means will put a cluster here.



Summary

- **Unsupervised learning**: fitting data without explicit labels.
- **Clustering**: finding 'groups' of related objects.
- **K-means**: simple iterative clustering strategy.
- **Vector quantization**: replacing measurements with 'prototypes'.
- **K-medians**: generalization to other distance functions.

- Next time:
 - Non-parametric clustering.