

UNSUPERVISED LEARNING

David Kauchak
CS 158 - Fall 2023

1

Administrative

Final project

- Project proposal feedback soon
- Progress report due next Tuesday

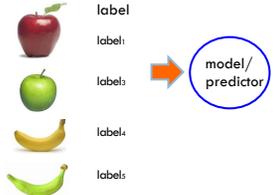
Mentor hours Thursday and Friday this week only

Monday office hours via zoom

No formal class Tuesday: working session for projects

2

Supervised learning



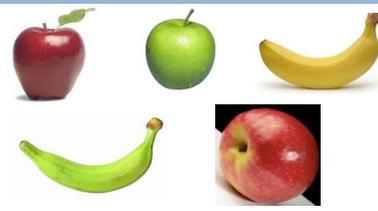
label
label
label
label
label

model/
predictor

Supervised learning: given labeled examples

3

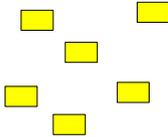
Unsupervised learning



Unsupervised learning: given data, i.e. examples, but no labels

4

Unsupervised learning



Given some example without labels, do something!

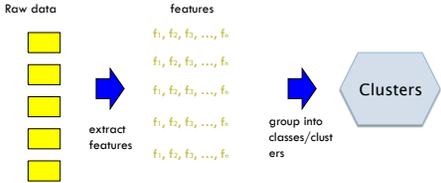
5

Unsupervised applications areas

- learn clusters/groups without any label
- customer segmentation (i.e. grouping)
- image compression
- bioinformatics: learn motifs
- find important features
- ...

6

Unsupervised learning: clustering



No "supervision", we're only given data and want to find natural groupings

7

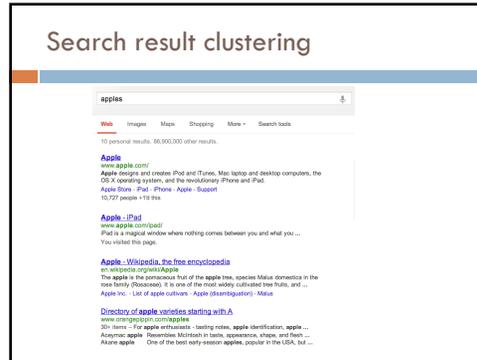
Unsupervised learning: modeling

Most frequently, when people think of unsupervised learning they think clustering

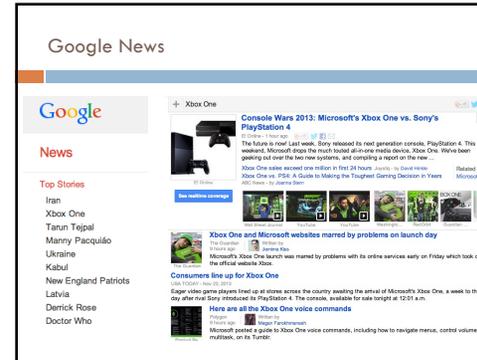
Another category: learning probabilities/parameters for models without supervision

- Learn a translation dictionary
- Learn a grammar for a language
- Learn the social graph

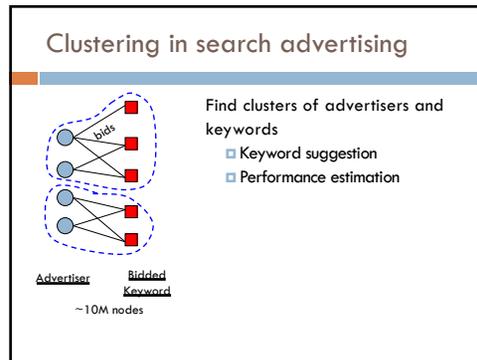
8



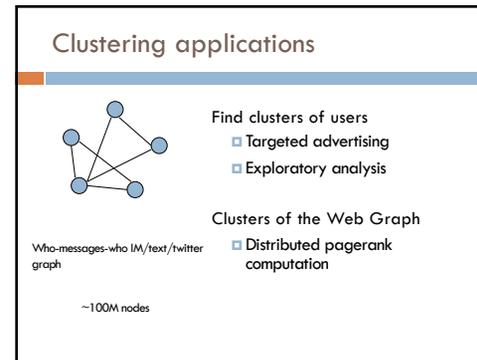
13



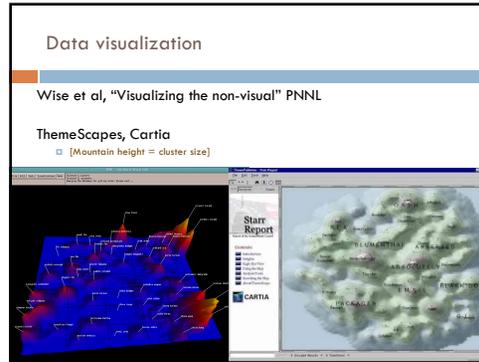
14



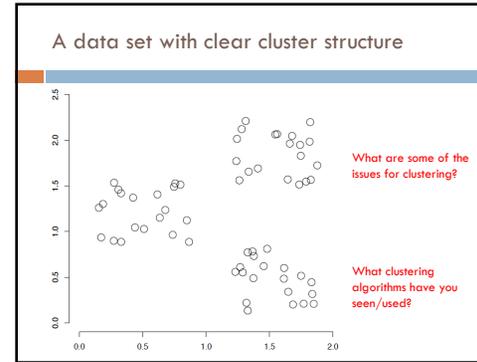
15



16



17



18

Issues for clustering

- Representation for clustering
 - How do we represent an example
 - features, etc.
 - Similarity/distance between examples
- Flat clustering or hierarchical
- Number of clusters
 - Fixed a priori
 - Data driven?

19

Clustering Algorithms

- Flat algorithms
 - Usually start with a random (partial) partitioning
 - Refine it iteratively
 - K means clustering
 - Model based clustering
 - Spectral clustering
- Hierarchical algorithms
 - Bottom-up, agglomerative
 - Top-down, divisive

The diagram shows two types of clustering visualizations. The top one shows three overlapping ellipses representing flat clustering. The bottom one shows a dendrogram representing hierarchical clustering.

20

Hard vs. soft clustering

Hard clustering: Each example belongs to exactly one cluster

Soft clustering: An example can belong to more than one cluster (probabilistic)

- Makes more sense for applications like creating browsable hierarchies
- You may want to put a pair of sneakers in two clusters: (i) sports apparel and (ii) shoes

21

K-means

Most well-known and popular clustering algorithm:

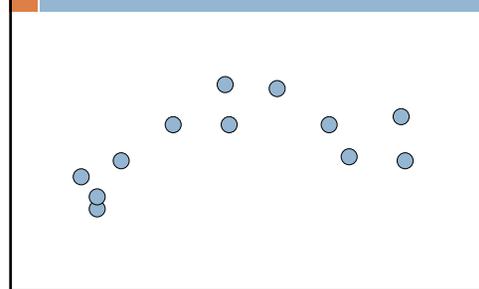
Start with some initial cluster centers

Iterate:

- Assign/cluster each example to closest center
- Recalculate centers as the mean of the points in a cluster

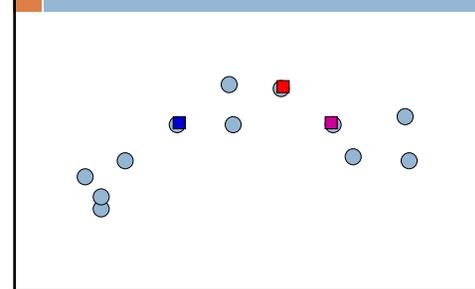
22

K-means: an example

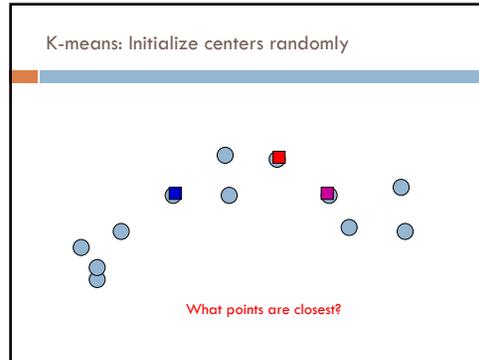


23

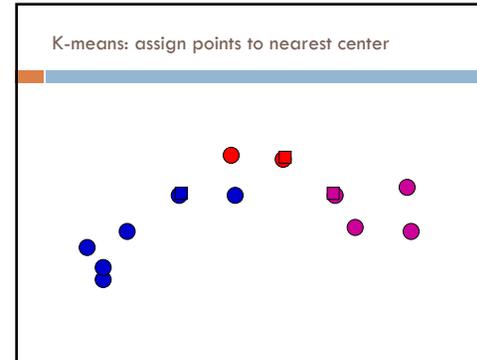
K-means: Initialize centers randomly



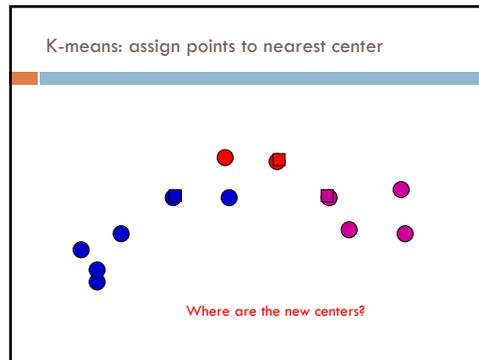
24



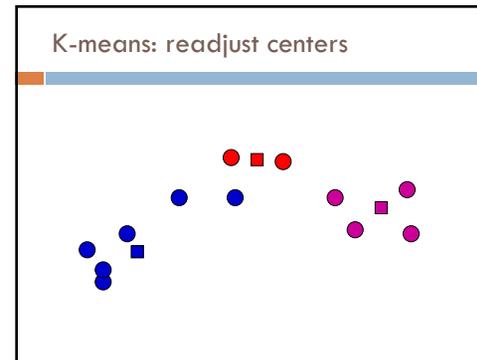
25



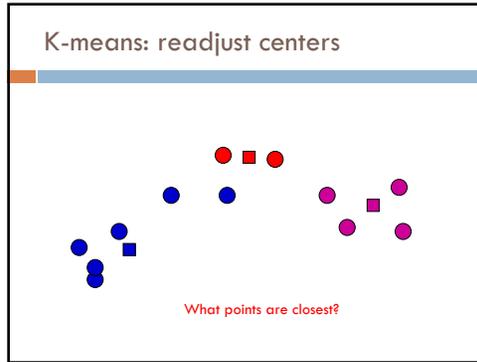
26



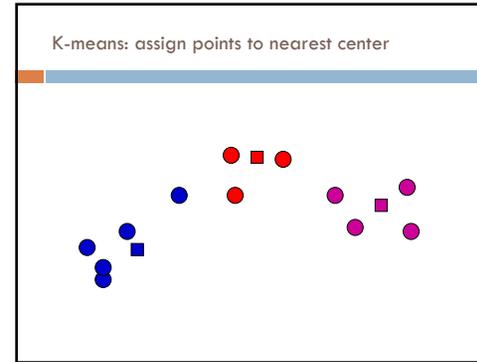
27



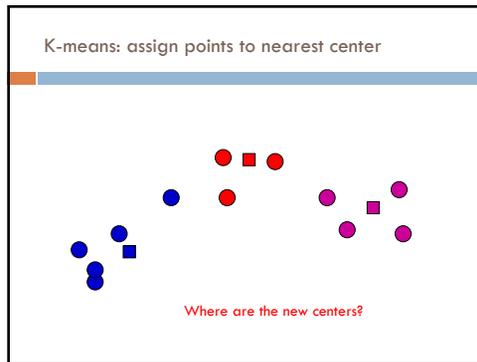
28



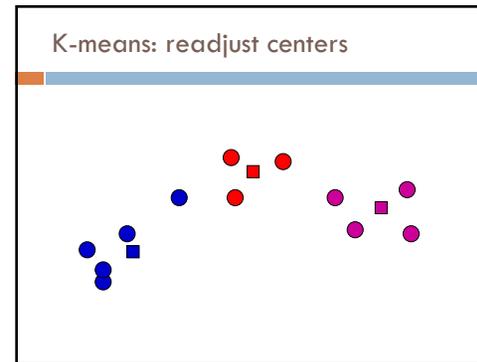
29



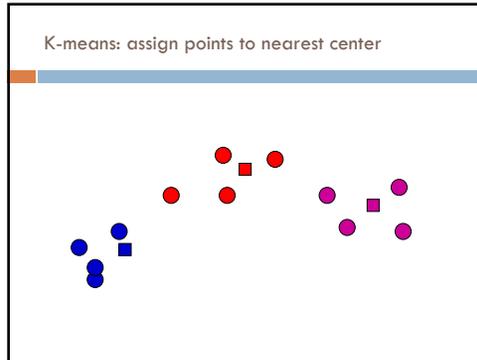
30



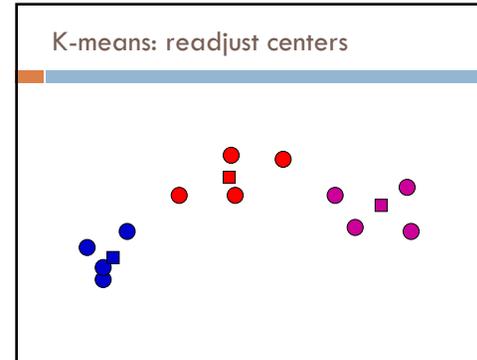
31



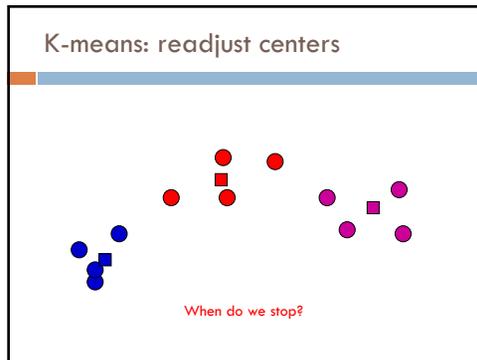
32



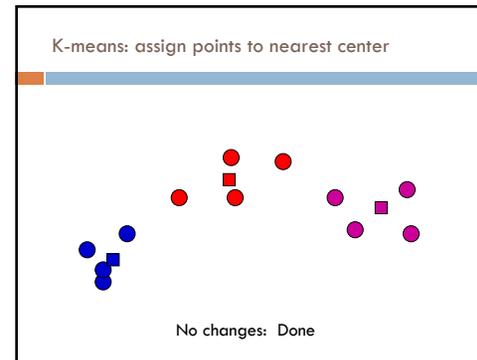
33



34



35



36

K-means

Iterate:

- Assign/cluster each example to closest center
- Recalculate centers as the mean of the points in a cluster

How do we do this?

37

K-means

Iterate:

- Assign/cluster each example to closest center
 - iterate over each point:
 - get distance to each cluster center
 - assign to closest center (hard cluster)
- Recalculate centers as the mean of the points in a cluster

38

K-means

Iterate:

- Assign/cluster each example to closest center
 - iterate over each point:
 - get distance to each cluster center
 - assign to closest center (hard cluster)
- Recalculate centers as the mean of the points in a cluster

What distance measure should we use?

39

Distance measures

Euclidean:

$$d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

good for spatial data

40

Clustering documents (e.g. wine data)

One feature for each word. The value is the number of times that word occurs.

Documents are points or vectors in this space

41

When Euclidean distance doesn't work

Which document is closest to q using Euclidean distance?

Which do you think should be closer?

42

Issues with Euclidean distance

the Euclidean distance between q and d_2 is large

but, the distribution of terms in q and d_2 are very similar

This is not what we want!

43

cosine similarity

$$\text{sim}(x,y) = \frac{x \cdot y}{|x||y|} = \frac{x \cdot y}{|x| \cdot |y|} = \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2} \sqrt{\sum_{i=1}^n y_i^2}}$$

correlated with the angle between two vectors

44

cosine distance

cosine similarity ranges from 0 and 1, with things that are similar 1 and dissimilar 0

cosine distance:

$$d(x, y) = 1 - \text{sim}(x, y)$$

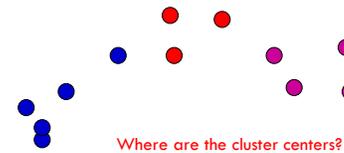
- good for text data and many other "real world" data sets
- *computationally friendly* since we only need to consider features that have non-zero values for **both** examples

45

K-means

Iterate:

- Assign/cluster each example to closest center
- Recalculate centers as the mean of the points in a cluster

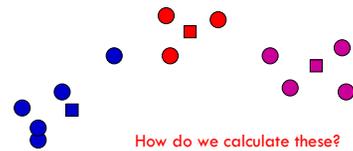


46

K-means

Iterate:

- Assign/cluster each example to closest center
- Recalculate centers as the mean of the points in a cluster



47

K-means

Iterate:

- Assign/cluster each example to closest center
- Recalculate centers as the mean of the points in a cluster

Mean of the points in the cluster:

$$\mu(C) = \frac{1}{|C|} \sum_{x \in C} x$$

where:

$$x + y = \sum_{i=1}^n x_i + y_i \quad \frac{x}{|C|} = \sum_{i=1}^n \frac{x_i}{|C|}$$

48

K-means loss function

K-means tries to minimize what is called the "k-means" loss function:

$$\text{loss} = \sum_{i=1}^n d(x_i, \mu_k)^2 \text{ where } \mu_k \text{ is cluster center for } x_i$$

the sum of the squared distances from each point to the associated cluster center

49

Minimizing k-means loss

Iterate:

1. Assign/cluster each example to closest center
2. Recalculate centers as the mean of the points in a cluster

$$\text{loss} = \sum_{i=1}^n d(x_i, \mu_k)^2 \text{ where } \mu_k \text{ is cluster center for } x_i$$

Does each step of k-means move towards reducing this loss function (or at least not increasing it)?

50

Minimizing k-means loss

Iterate:

1. Assign/cluster each example to closest center
2. Recalculate centers as the mean of the points in a cluster

$$\text{loss} = \sum_{i=1}^n d(x_i, \mu_k)^2 \text{ where } \mu_k \text{ is cluster center for } x_i$$

This isn't quite a complete proof/argument, but:

1. Any other assignment would end up in a larger loss
2. The mean of a set of values minimizes the squared error

51

Minimizing k-means loss

Iterate:

1. Assign/cluster each example to closest center
2. Recalculate centers as the mean of the points in a cluster

$$\text{loss} = \sum_{i=1}^n d(x_i, \mu_k)^2 \text{ where } \mu_k \text{ is cluster center for } x_i$$

Does this mean that k-means will always find the minimum loss/clustering?

52

Minimizing k-means loss

Iterate:

1. Assign/cluster each example to closest center
2. Recalculate centers as the mean of the points in a cluster

$$\text{loss} = \sum_{i=1}^n d(x_i, \mu_i)^2 \text{ where } \mu_i \text{ is cluster center for } x_i$$

NO! It will find a *minimum*.

Unfortunately, the k-means loss function is generally not convex and for most problems has many, many minima

We're only guaranteed to find one of them

53

K-means variations/parameters

Start with some initial cluster centers

Iterate:

- Assign/cluster each example to closest center
- Recalculate centers as the mean of the points in a cluster

What are some other variations/parameters we haven't specified?

54

K-means variations/parameters

Initial (seed) cluster centers

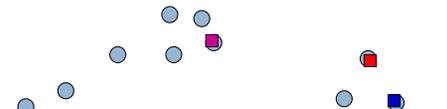
Convergence

- A fixed number of iterations
- partitions unchanged
- Cluster centers don't change

K!

55

K-means: Initialize centers randomly



What would happen here?

Seed selection ideas?

56

Seed choice

Results can vary drastically based on random seed selection

Some seeds can result in poor convergence rate, or convergence to sub-optimal clusterings

Common heuristics

- Random points (not examples) in the space
- Randomly pick examples
- Points least similar to any existing center (furthest centers heuristic)
- **Try out multiple starting points**
- Initialize with the results of another clustering method

57

Furthest centers heuristic

$\mu_1 =$ pick random point

for $i = 2$ to K :

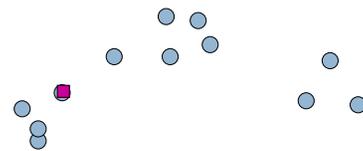
$\mu_i =$ point that is furthest from **any** previous centers

$$\mu_i = \underset{x}{\operatorname{arg\,max}} \min_{\mu_j : 1 < j < i} d(x, \mu_j)$$

point with the largest distance to any previous center
smallest distance from x to any previous center

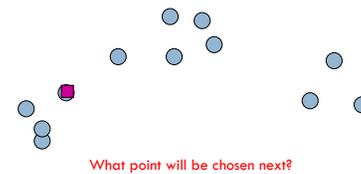
58

K-means: Initialize furthest from centers

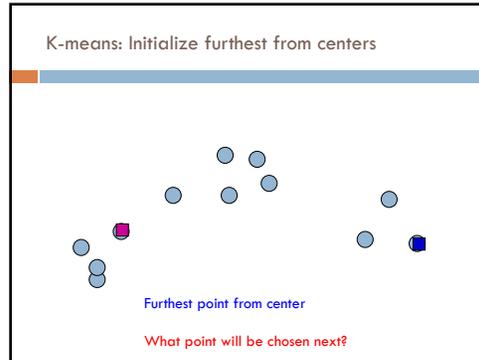


59

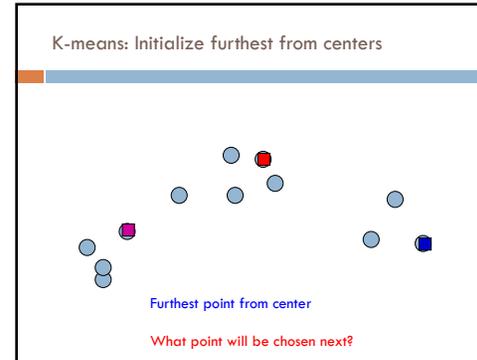
K-means: Initialize furthest from centers



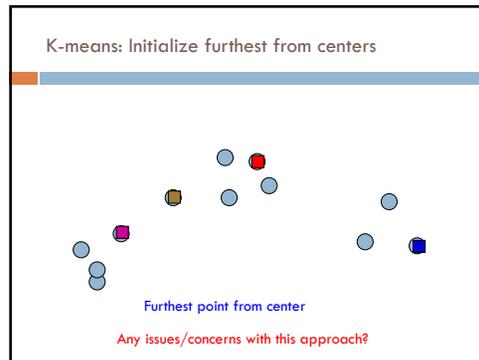
60



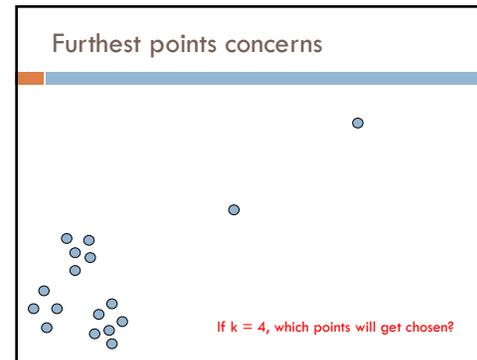
61



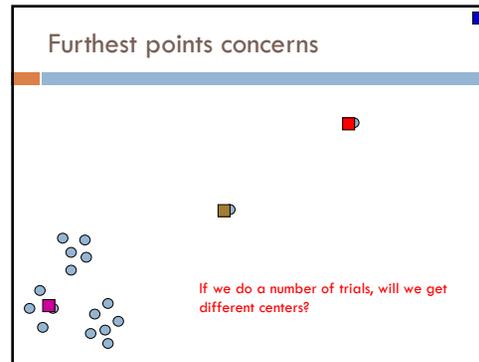
62



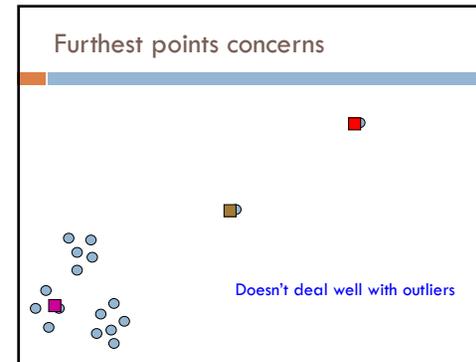
63



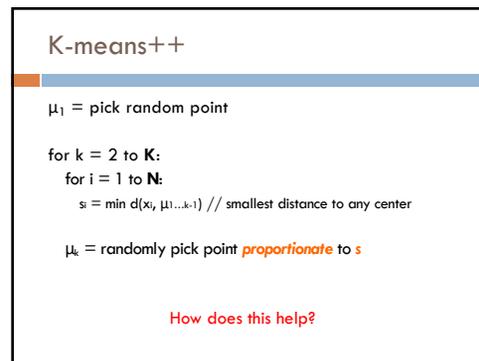
64



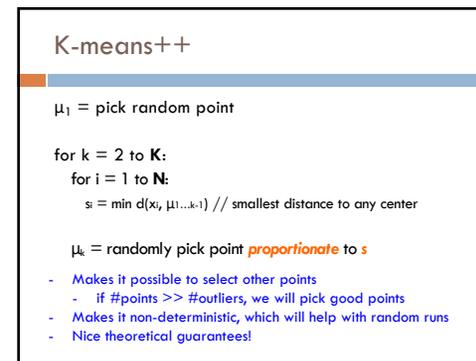
65



66



67



68

K-means variations/parameters

Initial (seed) cluster centers

Convergence

- ▣ A fixed number of iterations
- ▣ partitions unchanged
- ▣ Cluster centers don't change

K!

69

How Many Clusters?

Number of clusters K must be provided

How should we determine the number of clusters?

How did we deal with models becoming too complicated previously?

70

Many approaches

Regularization!!!

Statistical test

71

k-means loss revisited

K-means is trying to minimize:

$$\text{loss} = \sum_{i=1}^n d(x_i, \mu_i)^2 \text{ where } \mu_i \text{ is cluster center for } x_i$$

What happens when k increases?

72

k-means loss revisited

K-means is trying to minimize:

$$loss = \sum_{i=1}^n d(x_i, \mu_k)^2 \text{ where } \mu_k \text{ is cluster center for } x_i$$

Loss goes down!

Making the model more complicated allows us more flexibility, but can "overfit" to the data

73

k-means loss revisited

K-means is trying to minimize:

$$loss_{kmeans} = \sum_{i=1}^n d(x_i, \mu_k)^2 \text{ where } \mu_k \text{ is cluster center for } x_i$$



2 regularization options

$$loss_{BIC} = loss_{kmeans} + K \log N \quad (\text{where } N = \text{number of points})$$

$$loss_{AIC} = loss_{kmeans} + KN$$

What effect will this have?
Which will tend to produce smaller k?

74

k-means loss revisited

2 regularization options

$$loss_{BIC} = loss_{kmeans} + K \log N \quad (\text{where } N = \text{number of points})$$

$$loss_{AIC} = loss_{kmeans} + KN$$

AIC penalizes increases in K more harshly

Both require a change to the K-means algorithm

Tend to work reasonably well in practice if you don't know K

75