

Final project
Project proposal feedback on Gradescope
Progress report due Sunday

Midterm 2

No office hours next Mon/Tue

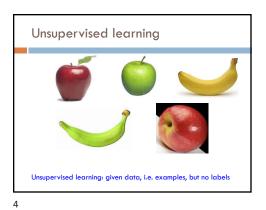
Class next Tuesday?

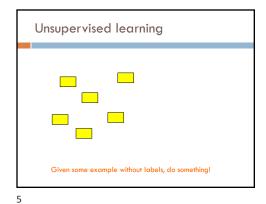
2

1

Supervised learning

label
labels
labels
Supervised learning: given labeled examples





Unsupervised applications areas

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

6

8

Unsupervised learning: clustering

Raw data

features

fi, f2, f3, ..., f6

fi, f2, f3, ..., f6

fi, f2, f3, ..., f6

group into classes/clust
extract
features

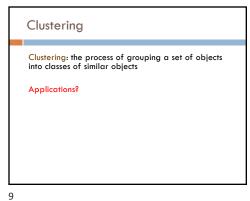
No "supervision", we're only given data and want to find natural groupings

Unsupervised learning: modeling

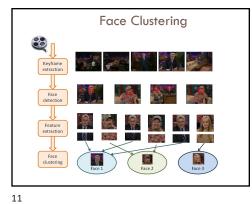
Most frequently, when people think of unsupervised learning they think of clustering

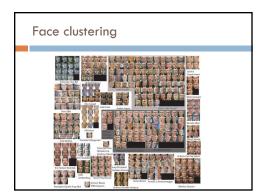
Another category: learning probabilities/parameters for models without supervision

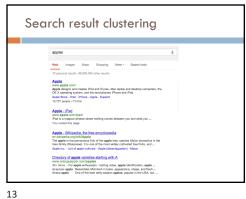
Learn a translation dictionary
Learn a grammar for a language
Learn the social graph



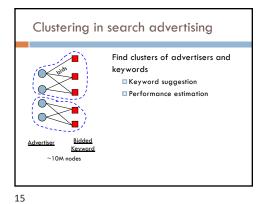


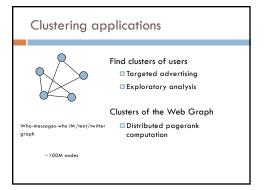


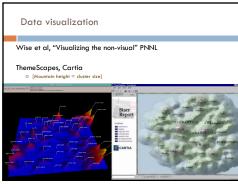












A data set with clear cluster structure What are some of the issues for clustering? 1.0 What clustering algorithms have you seen/used? 0.5 1.5

18

20

17

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



Hard clustering: Each example belongs to exactly one cluster

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

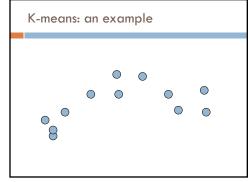
- Makes more sense for applications like creating browsable hierarchies
- You may want to put a pair of sneakers in two clusters: (i) sports apparel and (ii) shoes

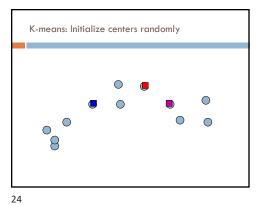
K-means Most well-known and popular clustering algorithm: Start with some initial cluster centers Iterate:

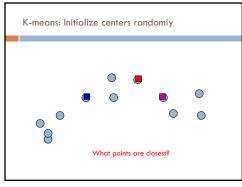
- Assign/cluster each example to closest center
- Recalculate centers as the mean of the points in a cluster

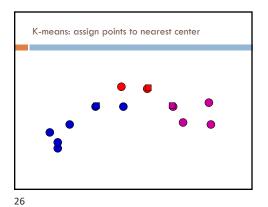
21

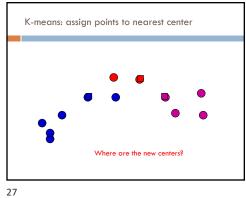
22

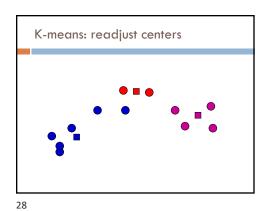


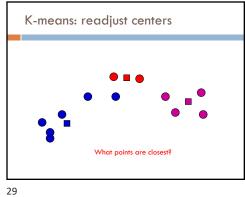


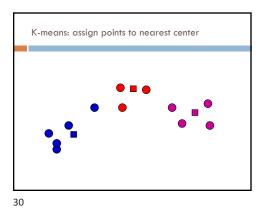


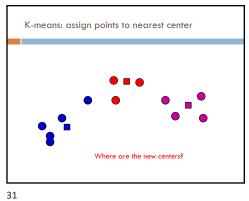


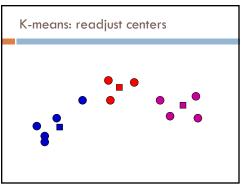


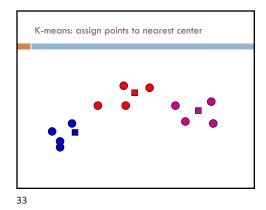


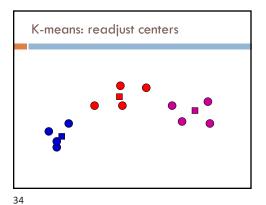






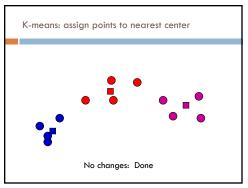




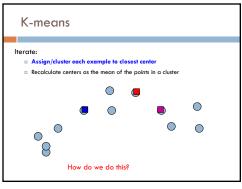


K-means: readjust centers

When do we stop?



35 36



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

40

37

39

Iterate:

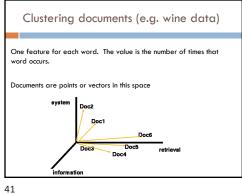
• Assign/cluster each example to closest center
iterate over each point:

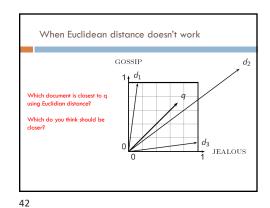
- get distance to each cluster center
- assign to closest center (hard cluster)

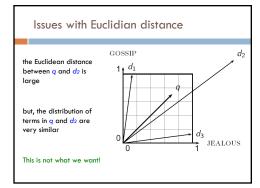
• Recalculate centers as the mean of the points in a cluster

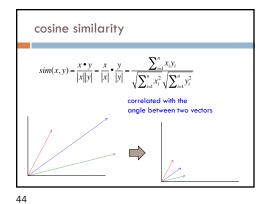
What distance measure should we use?

Distance measures $d(x,y) = \sqrt{\sum_{i=1}^{n} (x_i - y_i)^2}$ good for spatial data









cosine distance

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

cosine distance:

$$d(x,y) = 1 - sim(x,y)$$

- 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

K-means Iterate: Assign/cluster each example to closest center ☐ Recalculate centers as the mean of the points in a cluster Where are the cluster centers?

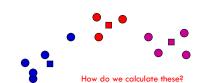
45

K-means

lterate:

□ Assign/cluster each example to closest center

Recalculate centers as the mean of the points in a cluster



47

K-means

Iterate:

46

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

K-means loss function

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

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

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

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

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

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

49

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

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

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

Minimizing k-means loss

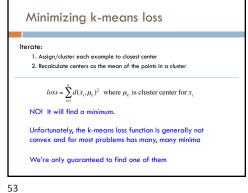
Iterate

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

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

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

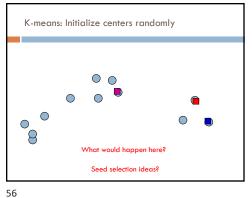
51

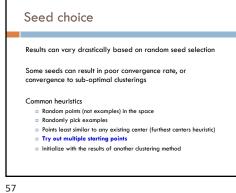


K-means variations/parameters Start with some initial cluster centers Iterate: Assign/cluster each example to closest center • Recalculate centers as the mean of the points in a cluster What are some other variations/parameters we haven't specified?

54

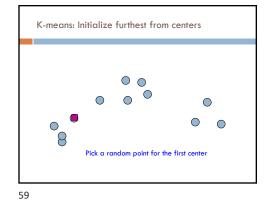
K-means variations/parameters Initial (seed) cluster centers Convergence ■ A fixed number of iterations partitions unchanged □ Cluster centers don't change K!

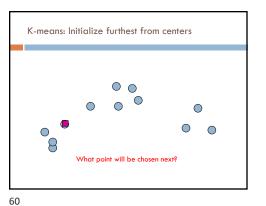


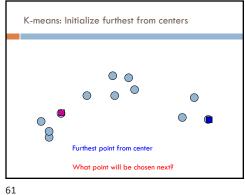


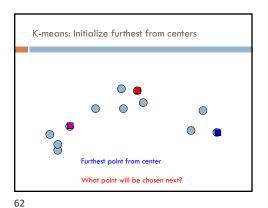
Furthest centers heuristic μ_1 = pick random point for i = 2 to K: $\mu_i = point \ that \ is \ furthest \ from \ \textbf{any} \ previous \ centers$ $\mu_i = \begin{array}{cc} \arg\max & \min \\ x & \mu_j : 1 < j < i \end{array} d(x, \mu_j)$ point with the largest distance smallest distance from x to any previous center to any previous center

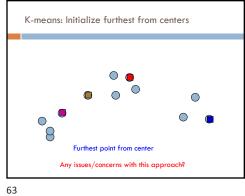
58

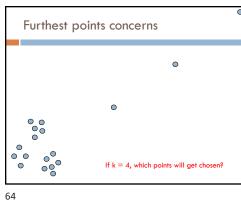


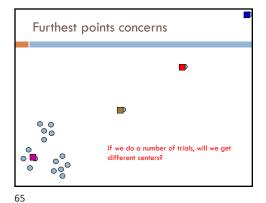


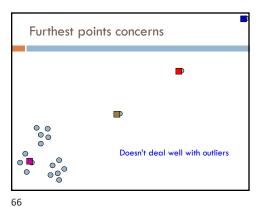












K-means++ $\mu_1 = \text{pick random point}$ for k = 2 to K:
for i = 1 to N: $s_i = \min d(x_i, \mu_{1...k-1}) // \text{smallest distance to any center}$ $\mu_k = \text{randomly pick point } \text{proportionate to s}$ How does this help?

K-means++

μ₁ = pick random point

for k = 2 to K:
 for i = 1 to N:
 si = min d(xi, μ1...λ1) // smallest distance to any center

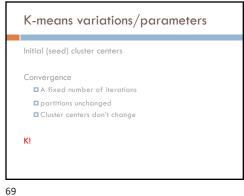
μ_k = randomly pick point proportionate to s

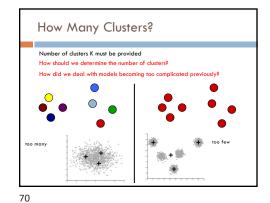
- Makes it possible to select other points

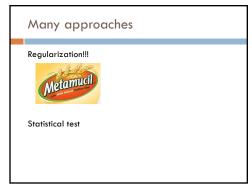
- if #points >> #outliers, we will pick good points

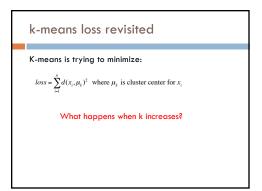
- Makes it non-deterministic, which will help with random runs

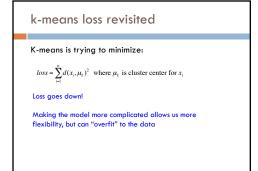
- Nice theoretical guarantees!











K-means loss revisited

K-means is trying to minimize: $loss_{lmeans} = \sum_{i=1}^{n} d(x_i, \mu_k)^2$ where μ_k is cluster center for x_i $loss_{BIC} = loss_{lmeans} + K \log N$ (where N = number of points) $loss_{AIC} = loss_{kmeans} + KN$ What effect will this have?

Which will tend to produce smaller k?

74

73

k-means loss revisited $loss_{BIC} = loss_{kmeans} + K \log N \qquad \text{(where N = number of points)}$ $loss_{AIC} = loss_{kmeans} + KN$ 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