# CLUSTERING BEYOND K-MEANS

David Kauchak
CS 158 – Spring 2022

1

## Administrative

Assignment 8 back

Final project status reports due Wednesday

Next class: skim the papers

2

## K-means

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

3

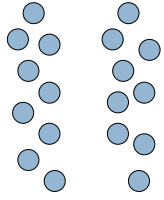## Problems with K-means

Determining K is challenging

Hard clustering isn't always right

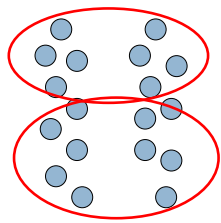Assumes clusters are spherical

Greedy approach

4

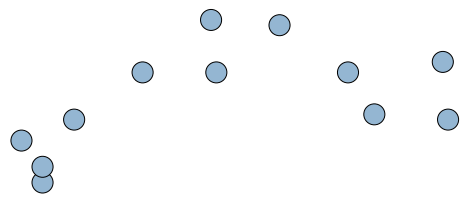## Problems with K-means

What would K-means give us here?

5

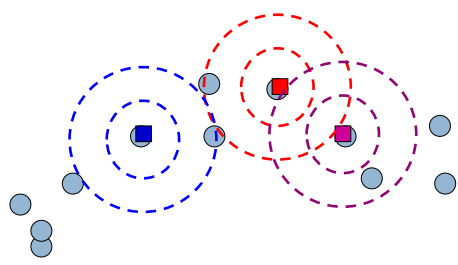## Assumes spherical clusters

k-means assumes spherical clusters!

6

## K-means: another view
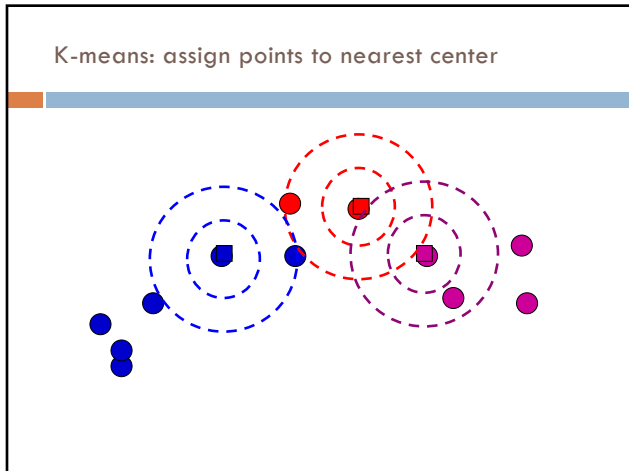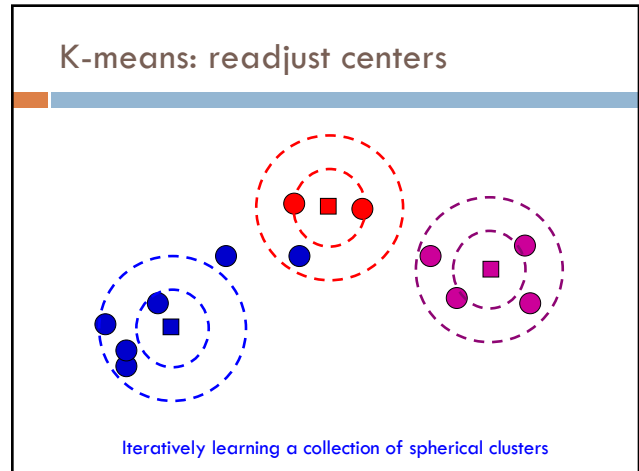
7

## K-means: another view
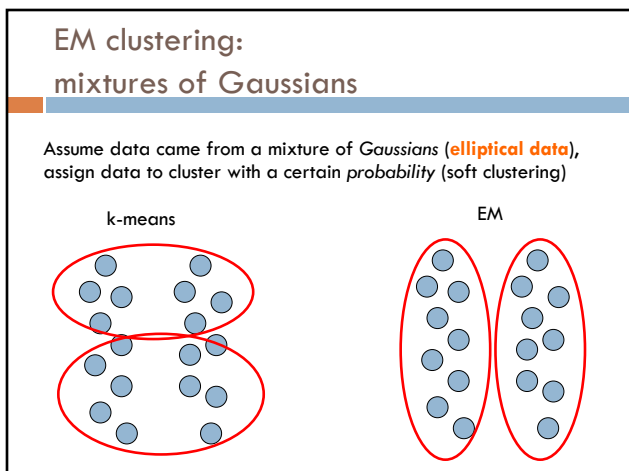
8

## K-means: assign points to nearest center



9

## K-means: readjust centers



Iteratively learning a collection of spherical clusters

10

## EM clustering:
## mixtures of Gaussians

Assume data came from a mixture of *Gaussians* (**elliptical data**), assign data to cluster with a certain *probability* (soft clustering)

k-means                                          EM



11

## EM clustering

Very similar at a high-level to K-means

Iterate between assigning points and recalculating cluster centers

Two main differences between K-means and EM clustering:
1. We assume elliptical clusters (instead of spherical)
2. It is a "soft" clustering algorithm

12

## Soft clustering



p(red) = 0.8
p(blue) = 0.2

p(red) = 0.9
p(blue) = 0.1

13

## EM clustering

Start with some initial cluster centers

*Iterate:*

- **soft assign** points to each cluster

    Calculate: $p(x; \theta_c)$

    the probability of each point belonging to each cluster

- recalculate the cluster centers

    Calculate new cluster parameters, $\theta_c$
    maximum likelihood cluster centers given the current soft clustering

14

## EM example



Start with some initial cluster centers

Figure from Chris Bishop

15

## Step 1: soft cluster points



Which points belong to which clusters (soft)?

Figure from Chris Bishop
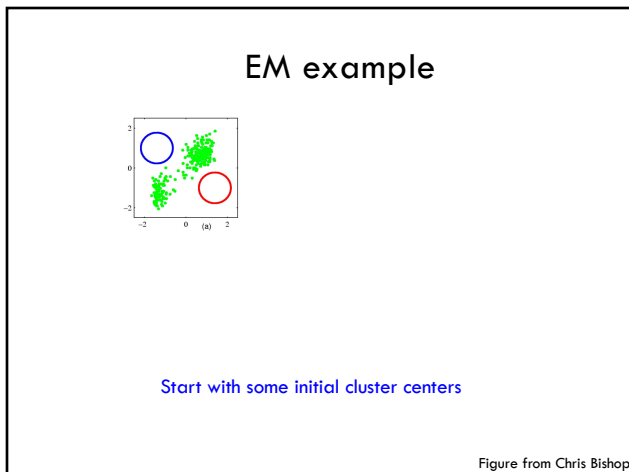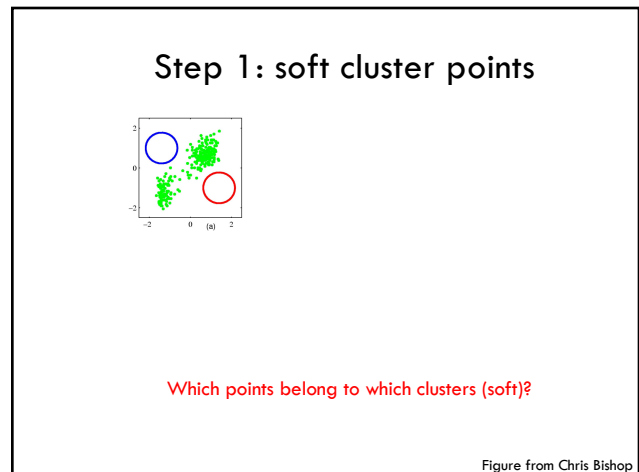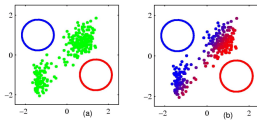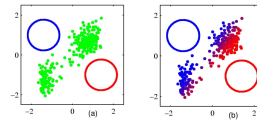
16

## Step 1: soft cluster points



Notice it's a soft (probabilistic) assignment

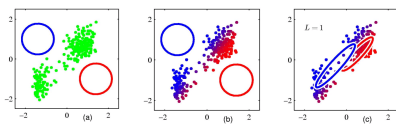Figure from Chris Bishop

17

## Step 2: recalculate centers



What do the new centers look like?

Figure from Chris Bishop

18

## Step 2: recalculate centers



Cluster centers get a **weighted** contribution from points

Figure from Chris Bishop
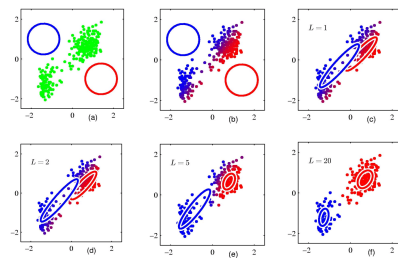
19

## keep iterating…



Figure from Chris Bishop
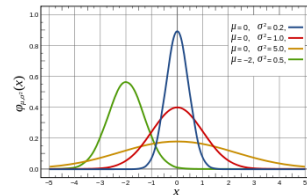
20

## Model: mixture of Gaussians

How do you define a Gaussian (i.e. ellipse)?
In 1-D?
In m-D?

21

## Gaussian in 1D

$$f(x; \sigma, \theta) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

parameterized by the mean and the standard deviation/variance

22

## Gaussian in multiple dimensions

$$N[x; \mu, \Sigma] = \frac{1}{(2\pi)^{d/2} \sqrt{\det(\Sigma)}} \exp[-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)]$$

Covariance determines
the shape of these contours

We learn the means of each cluster (i.e. the center) and the
covariance matrix (i.e. how spread out it is in any given direction)

23

## Step 1: soft cluster points

- soft assign points to each cluster
  Calculate: p(x; $\theta_c$)
  the probability of each point belonging to each cluster

  How do we calculate these probabilities?

24

## Step 1: soft cluster points



- soft assign points to each cluster
Calculate: p(x; $\theta_c$)
the probability of each point belonging to each cluster

Just plug into the Gaussian equation for each cluster!
(and normalize to make a probability)

25

## Step 2: recalculate centers



Recalculate centers:
calculate new cluster parameters, $\theta_c$
maximum likelihood cluster centers given the current soft clustering

How do calculate the cluster centers?

26

## Fitting a Gaussian

What is the "best"-fit Gaussian for this data?

10, 10, 10, 9, 9, 8, 11, 7, 6, …

Recall this is the 1-D Gaussian equation:

$$f(x;\sigma,\theta) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

27

## Fitting a Gaussian

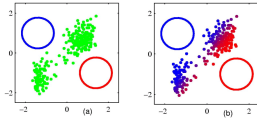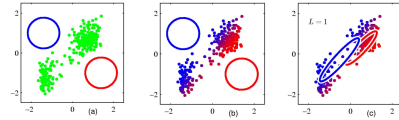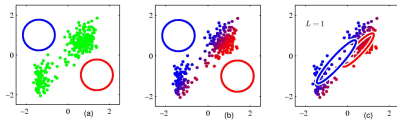What is the "best"-fit Gaussian for this data?

10, 10, 10, 9, 9, 8, 11, 7, 6, …

The MLE is just the mean and variance of the data!

Recall this is the 1-D Gaussian equation:

$$f(x;\sigma,\theta) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$
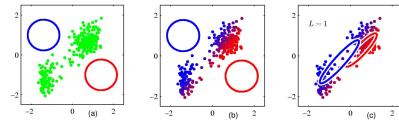
28

## Step 2: recalculate centers



Recalculate centers:
Calculate $\theta_c$
maximum likelihood cluster centers given the current
*soft clustering*

How do we deal with "soft" data points?

29

## Step 2: recalculate centers



Recalculate centers:
Calculate $\theta_c$
maximum likelihood cluster centers given the current
*soft clustering*

Use fractional counts!

30

## E and M steps: creating a better model

EM stands for Expectation Maximization

**Expectation**: Given the current model, figure out the expected probabilities of the data points to each cluster

$p(x; \theta_c)$ What is the probability of each point belonging to each cluster?

**Maximization**: Given the probabilistic assignment of all the points, estimate a new model, $\theta_c$

Just like NB maximum likelihood estimation, except we use fractional counts instead of whole counts

31

## Similar to *k*-means

Iterate:
Assign/cluster each point to closest center

Expectation: Given the current model, figure out the expected probabilities of the points to each cluster $p(x; \theta_c)$

Recalculate centers as the mean of the points in a cluster

Maximization: Given the probabilistic assignment of all the points, estimate a new model, $\theta_c$

32

8

## E and M steps

**Expectation**: Given the current model, figure out the expected probabilities of the data points to each cluster

**Maximization**: Given the probabilistic assignment of all the points, estimate a new model, $\theta_c$

*Iterate:*

each iterations increases the likelihood of the data and is guaranteed to converge (though to a local optimum)!

33

## EM

EM is a general purpose approach for training a model when you don't have labels

Not just for clustering!
- K-means is just for clustering

One of the most general purpose unsupervised approaches
- can be hard to get right!

34

## EM is a general framework

Create an initial model, $\theta'$
- Arbitrarily, randomly, or with a small set of training examples

Use the model $\theta'$ to obtain another model $\theta$ such that

$\sum_i \log P_\theta(data_i) > \sum_i \log P_{\theta'}(data_i)$     i.e. better models data (increased log likelihood)

Let $\theta' = \theta$ and repeat the above step until reaching a local maximum
- Guaranteed to find a better model after each iteration

Where else have you seen EM?

35

## EM shows up all over the place

Training HMMs (Baum-Welch algorithm)

Learning probabilities for Bayesian networks

EM-clustering

Learning word alignments for language translation

Learning Twitter friend network

Genetics

Finance

Anytime you have a model and unlabeled data!

36

## Finding Word Alignments

… la maison … la maison bleue … la fleur …

… the house … the blue house … the flower …

In machine translation, we train from pairs of translated sentences

Often useful to know how the words align in the sentences

Use EM!
• learn a model of P(french-word | english-word)

37

## Finding Word Alignments

… la maison … la maison bleue … la fleur …



… the house … the blue house … the flower …

All word alignments are equally likely

All P(french-word | english-word) equally likely

38

## Finding Word Alignments

… la maison … la maison bleue … la fleur …



… the house … the blue house … the flower …

"la" and "the" observed to co-occur frequently, so P(la | the) is increased.

39

## Finding Word Alignments

… la maison … la maison bleue … la fleur …
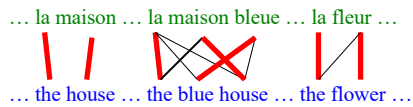


… the house … the blue house … the flower …

"house" co-occurs with both "la" and "maison", but P(maison | house) can be raised without limit, to 1.0, while P(la | house) is limited because of "the"

(pigeonhole principle)

40

## Finding Word Alignments

… la maison … la maison bleue … la fleur …

… the house … the blue house … the flower …

settling down after another iteration

41

## Finding Word Alignments

… la maison … la maison bleue … la fleur …

… the house … the blue house … the flower …

**Inherent hidden structure revealed by EM training!**
For details, see
- "A Statistical MT Tutorial Workbook" (Knight, 1999).
  - 37 easy sections, final section promises a free beer.
- "The Mathematics of Statistical Machine Translation"
  (Brown et al, 1993)
- Software:  GIZA++

42

## Statistical Machine Translation

… la maison … la maison bleue … la fleur …

… the house … the blue house … the flower …

P(maison | house ) = 0.411
P(maison | building) = 0.027
P(maison | manson) = 0.020
…

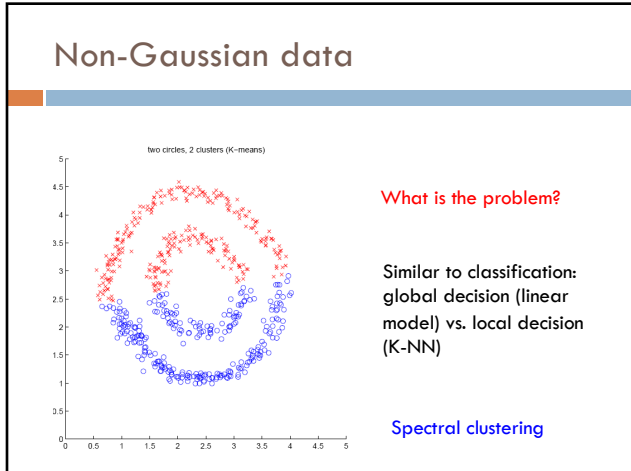Estimating the model from training data

43

## Other clustering algorithms

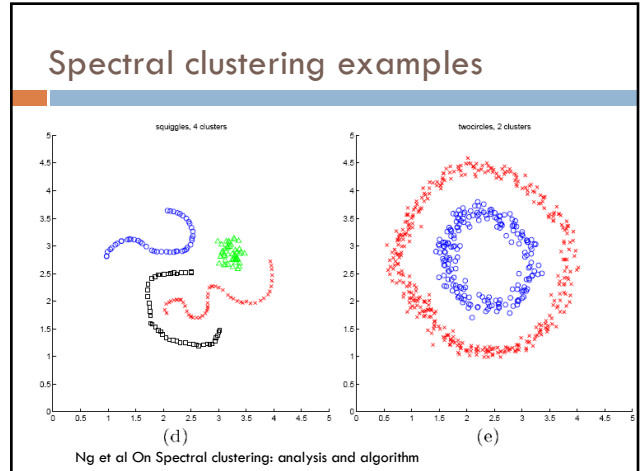K-means and EM-clustering are by far the most popular for clustering

However, they can't handle all clustering tasks

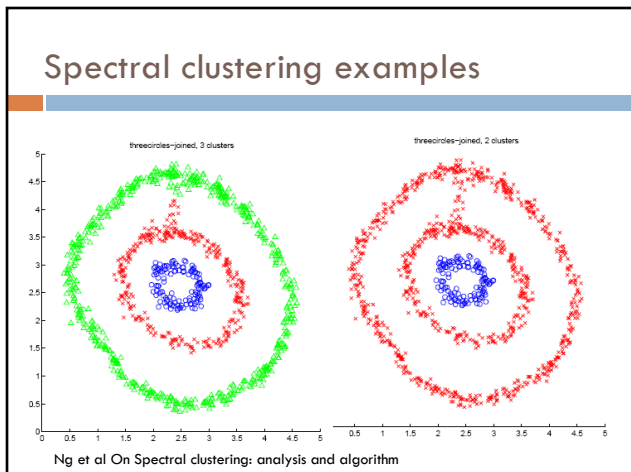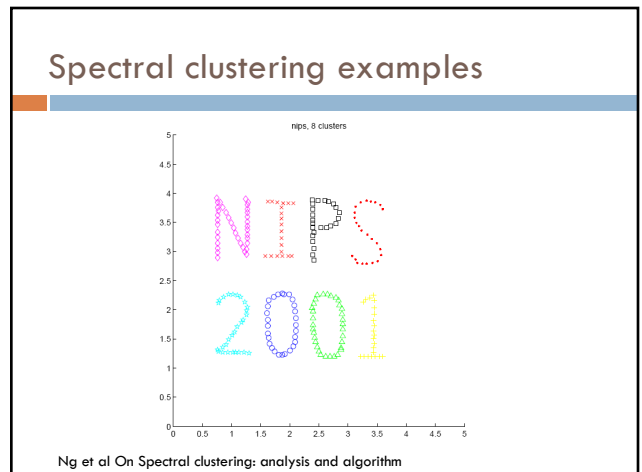What types of clustering problems can't they handle?

44

## Non-Gaussian data



two circles, 2 clusters (K-means)

**What is the problem?**

Similar to classification: global decision (linear model) vs. local decision (K-NN)

**Spectral clustering**

45

## Spectral clustering examples



squiggles, 4 clusters        twocircles, 2 clusters

(d)                          (e)

Ng et al On Spectral clustering: analysis and algorithm

46

## Spectral clustering examples



threecircles–joined, 3 clusters        threecircles–joined, 2 clusters

Ng et al On Spectral clustering: analysis and algorithm

47

## Spectral clustering examples



nips, 8 clusters

Ng et al On Spectral clustering: analysis and algorithm

48

## What Is A Good Clustering?

Internal criterion: A good clustering will produce high quality clusters in which:

- the <u>intra-class</u> (that is, intra-cluster) similarity is high
- the <u>inter-class</u> similarity is low
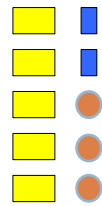
How would you evaluate clustering?

49

## Common approach: use labeled data

Use data with known classes

- For example, document classification data

data    label

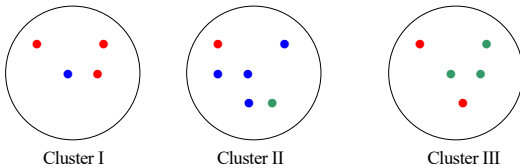If we clustered this data (ignoring labels) what would we like to see?

Reproduces class partitions

How can we quantify this?

50

## Common approach: use labeled data

**Purity**, the proportion of the dominant class in the cluster

Cluster I          Cluster II          Cluster III

Cluster I: Purity = (max(3, 1, 0)) / 4 = 3/4
Cluster II: Purity = (max(1, 4, 1)) / 6 = 4/6
Cluster III: Purity = (max(2, 0, 3)) / 5 = 3/5

Overall purity?

51

## Overall purity

Cluster I: Purity = (max(3, 1, 0)) / 4 = 3/4
Cluster II: Purity = (max(1, 4, 1)) / 6 = 4/6
Cluster III: Purity = (max(2, 0, 3)) / 5 = 3/5

Cluster average:

$$\frac{\frac{3}{4}+\frac{4}{6}+\frac{3}{5}}{3}=0.672$$

Weighted average: $\dfrac{4*\frac{3}{4}+6*\frac{4}{6}+5*\frac{3}{5}}{15}=\dfrac{3+4+3}{15}=0.667$
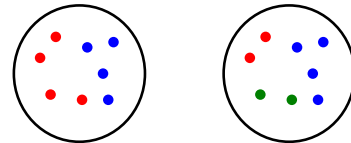
52

## Purity issues…

**Purity**, the proportion of the dominant class in the cluster

Good for comparing two algorithms, but not understanding how well a single algorithm is doing, why?

- Increasing the number of clusters increases purity

53

## Purity isn't perfect



Which is better based on purity?

Which do you think is better?

Ideas?

54

## Common approach: use labeled data

**Average entropy** of classes in clusters
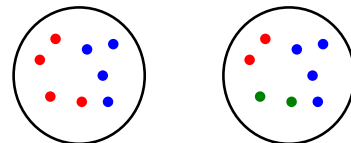
$$entropy(cluster) = -\sum_i p(class_i) \log p(class_i)$$

where p(class$_i$) is proportion of class $i$ in cluster

55

## Common approach: use labeled data

**Average entropy** of classes in clusters

$$entropy(cluster) = -\sum_i p(class_i) \log p(class_i)$$
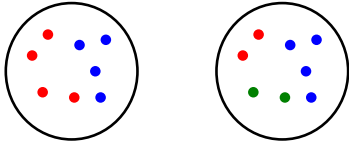


entropy?

56

## Common approach: use labeled data

**Average entropy** of classes in clusters

$$entropy(cluster) = -\sum_i p(class_i)\log p(class_i)$$

$-0.5\log 0.5 - 0.5\log 0.5 = 1$      $-0.5\log 0.5 - 0.25\log 0.25 - 0.25\log 0.25 = 1.5$

57

## Where we've been!

How many slides?

1,385 slides

58

## Where we've been!

Our ML suite:

How many classes?

How many lines of code?

59

## Where we've been!

Our ML suite:

29 classes

2951 lines of code

60

## Where we've been!

Our ML suite:
- Supports 7 classifiers
  - Decision Tree
  - Perceptron
  - Average Perceptron
  - Gradient descent
    - 2 loss functions
    - 2 regularization methods
  - K-NN
  - Naïve Bayes
  - 2 layer neural network
- Supports two types of data normalization
  - feature normalization
  - example normalization
- Supports two types of meta-classifiers
  - OVA
  - AVA

61

## Where we've been!

Hadoop!

- 532 lines of hadoop code in demos

62

## Where we've been!

Geometric view of data

Model analysis and interpretation (linear, etc.)

Evaluation and experimentation

Probability basics

Regularization (and priors)

Deep learning

Ensemble methods

Unsupervised learning (clustering)

63