Today's material available at:

https://github.com/juliennelachance/ai4all_clustering
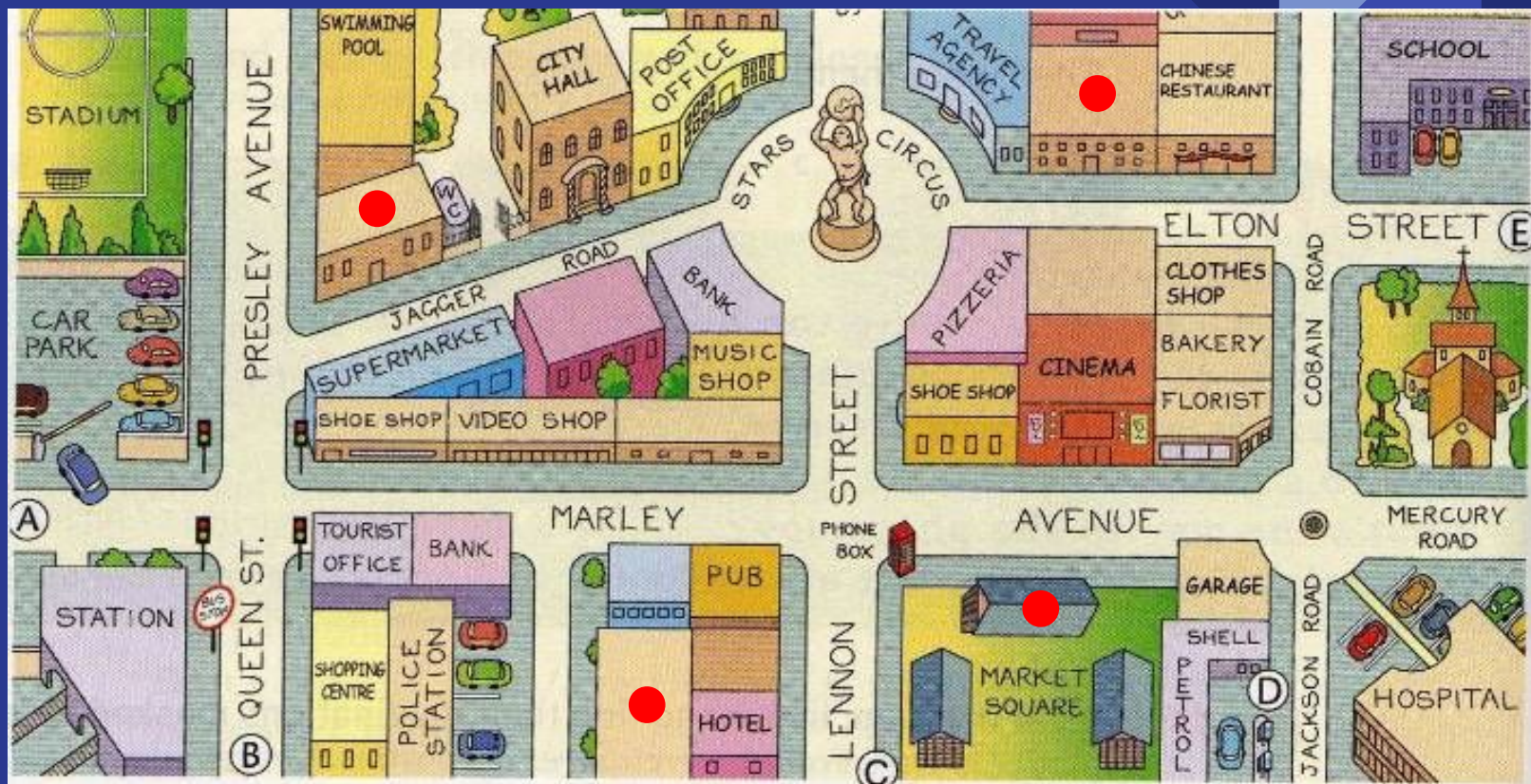
# Clustering

Princeton AI4ALL

By Julie LaChance

# Icebreaker

- Everyone receives a slip of paper

- Your task:
    Using the information about your person, form dog-walking groups

- Be prepared to justify your answers!
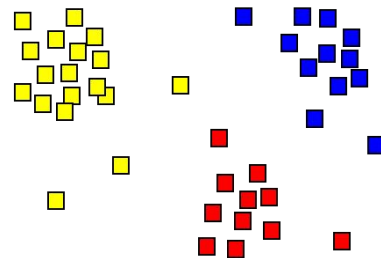    - What features are important and why?

# Overview:

- What is clustering?
- Applications of clustering
- Types of clustering
- k-means algorithm
- Recommendation system example

# What is clustering?

# What is clustering?



- From Wikipedia:

    "**Cluster analysis** or **clustering** is the task of grouping a set of objects in such a way that objects in the same group (called a **cluster**) are more similar (in some sense) to each other than to those in other groups (clusters)."

- Thinking back to our first lecture on machine learning:

    This is an important type of ***unsupervised learning***. Why?

# Applications of clustering

# Periodic Table of Elements

# Earthquakes and Seismology

# Medical Imaging

# Recommendation Systems: Amazon

# Recommendation Systems: Amazon

- Amazon doesn't know what it's like to read a book, or what you feel like when you read a particular book

- Amazon *does* know that people who bought a certain book also bought other books

- Patterns in the data can used to make recommendations

- If you've built up a long purchase history you'll often see pretty sophisticated recommendations

# And even more nefarious purposes...

# Discussion:

Can you come up with other applications of clustering?

# Types of clustering

# Two main types of clustering:

- **Hard clustering**:
    Each object is in exactly one class.

- **Soft clustering**:
    Objects are assigned a degree to which they are in each class.
    - (Think of this as probability of being in a class)

# Flat vs. Hierarchical clustering

- Flat = distinct clusters
- Hierarchical = clusters within clusters

More later on types of clustering *algorithms*. But first, an example...

# Campgrounds

# Campgrounds

# Campgrounds

# Campgrounds

# Breakout Session:

Try out the demo at
http://mlehman.github.io/kmeans-javascript/

Discuss with your table: how does the algorithm work?

(Note: this demo has a bug! Can't use more than 6 clusters)

# k-means clustering

# How can we automatically cluster campgrounds?



*k*-means clustering

Goal: Assign each of the *n* points to one out of *k* clusters (defined by a cluster centroid).

# Find 5 clusters for the 1000 points below!

# Step 0: Pick k random points as the cluster centroids

# Iteration 1



Step 1: Assign each point to closest cluster centroid

# Iteration 1 continued



Step 2: Find cluster centroids as center of mass of elements

# Iteration 2



Step 1: Assign each point to closest cluster centroid

# Iteration 2 continued



Step 2: Find cluster centroids as center of mass of elements

# Iteration 3



Step 1: Assign each point to closest cluster centroid

# Iteration 3 continued



Step 2: Find cluster centroids as center of mass of elements

# Iteration 4



Step 1: Assign each point to closest cluster centroid

No change in cluster assignment. TERMINATE!

# k-means overview

- Choose the value of k
  - This is how many cluster centroids we're finding
- Choose k points in the set
  - These are the initial centroid locations
- For each point not selected, assign to its nearest centroid
  - Now all points have an initial cluster assignment
- Until "happy" do:
  - 1. Recompute centroids of clusters
  - 2. Reassign all points to closest centroid (forms new clusters)

# Discussion:

Is this "hard" or "soft" clustering?

Is this "flat" or "hierarchical" clustering?

# Discussion:

When would k-means fail?
When would it succeed?

What are some drawbacks?
What are some advantages?

# Weaknesses:

## Non-globular clusters

# Weaknesses:

# Weaknesses:

Outliers and empty clusters

# Algorithm design: Distance metrics and features

- How do we determine what elements are "close"?
    - Spatial distance, similarity of campsite attributes,...

- "Features" are attributes that we use to mathematically compute closeness.
    - Height, age, hometown

- Weights for features
    - Are all features equally important or are some more important than others?
    - Are the features on different scales?

# Types of features

- **Categorical features** take a fixed set of values. The values cannot be ordered.
  - For example, the type of phone (android, iphone, windows phone).

- **Ordinal features** take a fixed set of values. The values can be ordered.
  - For example, the ranks of football teams in the PAC-12 conference.

- **Continuous features** can take any real value.
  - For example, the distance of Princeton University from your hometown.

# Main types of clustering algorithms

- Connectivity-based (hierarchical) clustering
- Centroid-based clustering
- Distribution-based clustering
- Density-based clustering

# Connectivity-based (hierarchical) clustering

- A hierarchy of clusters based on distance:



- Typically, form a "dendogram"

# Centroid-based clustering

- k-means!

# Distribution-based clustering

- Cluster according to how likely it is that points lie in a certain distribution

# Density-based clustering

- Group together points that are closely spaced, and marks farther-out points as outliers

outliers

# Recommendation system example

# Dataset: Movie critics

| Critic | Star Wars | Raiders of the Lost Arc | Casablanca | Singin' in the Rain |
|--------|-----------|-------------------------|------------|---------------------|
| Sam | **** | **** | * | ** |
| Sandy | ***** | **** | ** | * |
| Matt | ** | ** | **** | *** |
| Julia | ** | * | *** | **** |
| Sarah | ***** | ? | ? | ** |

- Could an algorithm use this data to recommend movies?
- How would you do it?

# Dataset: Movie critics

| Critic | Star Wars | Raiders of the Lost Arc | Casablanca | Singin' in the Rain |
|--------|-----------|-------------------------|------------|---------------------|
| Sam | **** | **** | * | ** |
| Sandy | ***** | **** | ** | * |
| Matt | ** | ** | **** | *** |
| Julia | ** | * | *** | **** |
| Sarah | ***** | ? | ? | ** |

- Could an algorithm use this data to recommend movies?
- How would you do it?

# Critics with similar taste

■ Preference space

# Measuring distance

**Star Wars**

5                          Sam
4                          Sandy
3
2              Matt
1      Julia
       1      2      3      4      5      **Raiders**

$$d_{ij} = \sqrt{(x_i - x_j)^2 + (y_i - y_j)^2}$$

$(x_2, y_2)$

$(x_1, y_1)$

- Measure similarity between points using a measure of distance

# Finding critics with similar taste

**Star Wars**

5            Sam
4            Sandy
3
2        Matt
1    Julia

  1    2    3    4    5   **Raiders**

$(x_2, y_2)$

$(x_1, y_1)$

$$d_{ij} = \sqrt{(x_i - x_j)^2 + (y_i - y_j)^2}$$

- People who liked Star Wars are close in preference space to those who liked Raiders

# Making a recommendation

**Star Wars**

5                 Sarah

                  Sam

4                 Sandy

3

2        Matt

1    Julia

      1     2     3     4     5   **Raiders**

$(x_2, y_2)$

$(x_1, y_1)$

$$d_{ij} = \sqrt{(x_i - x_j)^2 + (y_i - y_j)^2}$$

- Sarah hasn't seen Raiders, but gave Star Wars five stars.
- Chances are she'll like Raiders too!

# Features

- We used features to compare critics

- Feature: a data attribute used to make a comparison

- Quantify attributes of an object (size, weight, color, shape, density) in a way a computer can understand

- Quality is important

# Apples vs. Oranges



- A good feature discriminates between classes

- Think: how well does a feature help us tell two things apart?

- Is mass a good feature? By itself?

- What about in conjunction with another feature like color?

# Features to compare movies

| Feature | Star Wars | Raiders of the Lost Arc | Casablanca | Singin' in the Rain |
|---|---|---|---|---|
| … | | | | |
| … | | | | |
| … | | | | |
| … | | | | |
| … | | | | |

- Can you suggest some features to compare movies?

# Features to compare movies

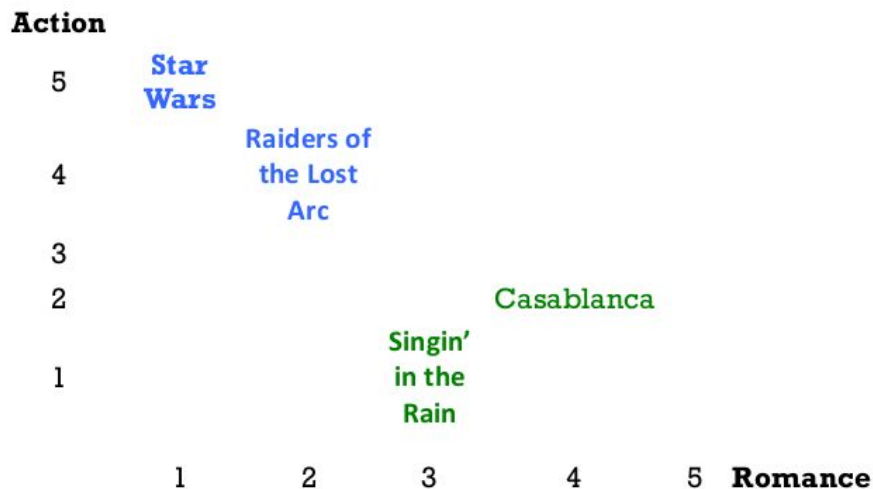| Feature | Star Wars | Raiders of the Lost Arc | Casablanca | Singin' in the Rain |
|---|---|---|---|---|
| Action (1 to 5) | 5 | 4 | 2 | 1 |
| Romance (1 to 5) | 1 | 2 | 4 | 3 |
| Length (min) | 121 | 115 | 102 | 103 |
| Harrison Ford | Y | Y | N | N |
| Year | 1977 | 1981 | 1942 | 1952 |

- What type of features are these?

# Comparing movies in feature space



- Here we can "eyeball" the clusters, but what if we had more features? (i.e. higher dimensional data)?

Exercise time:

Jupyter notebook available at
https://github.com/juliennelachance/ai4all_clustering