

Introduction to Machine Learning

Part 1 and Part 2

Yingyu Liang

`yliang@cs.wisc.edu`

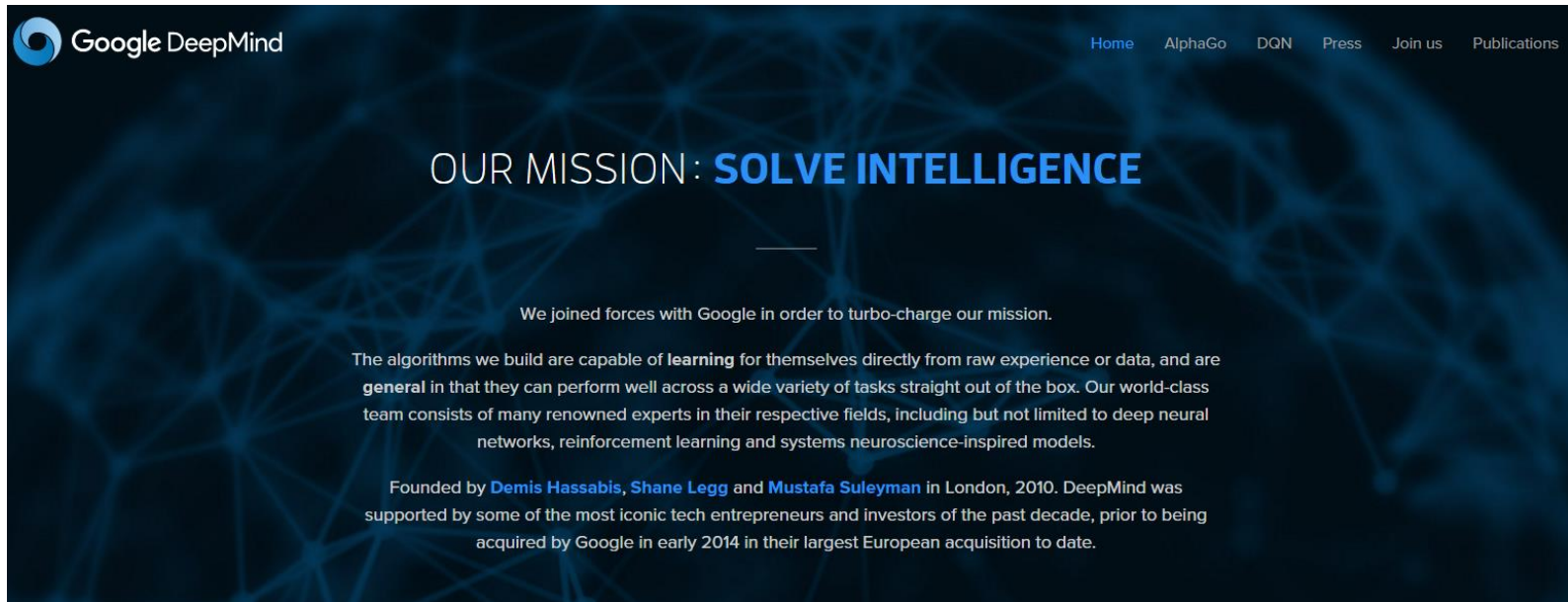
**Computer Sciences Department
University of Wisconsin, Madison**

What is machine learning?

- Short answer: recent buzz word

Industry

- Google

The image shows a screenshot of the Google DeepMind website. The background is dark blue with a network-like pattern of glowing lines. In the top left corner, there is the Google DeepMind logo. In the top right corner, there is a navigation menu with links for Home, AlphaGo, DQN, Press, Join us, and Publications. The main heading reads "OUR MISSION: SOLVE INTELLIGENCE". Below this, there is a sub-heading "We joined forces with Google in order to turbo-charge our mission." followed by a paragraph describing the company's mission and team. At the bottom, there is a paragraph about the company's founding and acquisition by Google.

Google DeepMind

Home AlphaGo DQN Press Join us Publications

OUR MISSION: **SOLVE INTELLIGENCE**

We joined forces with Google in order to turbo-charge our mission.

The algorithms we build are capable of **learning** for themselves directly from raw experience or data, and are **general** in that they can perform well across a wide variety of tasks straight out of the box. Our world-class team consists of many renowned experts in their respective fields, including but not limited to deep neural networks, reinforcement learning and systems neuroscience-inspired models.

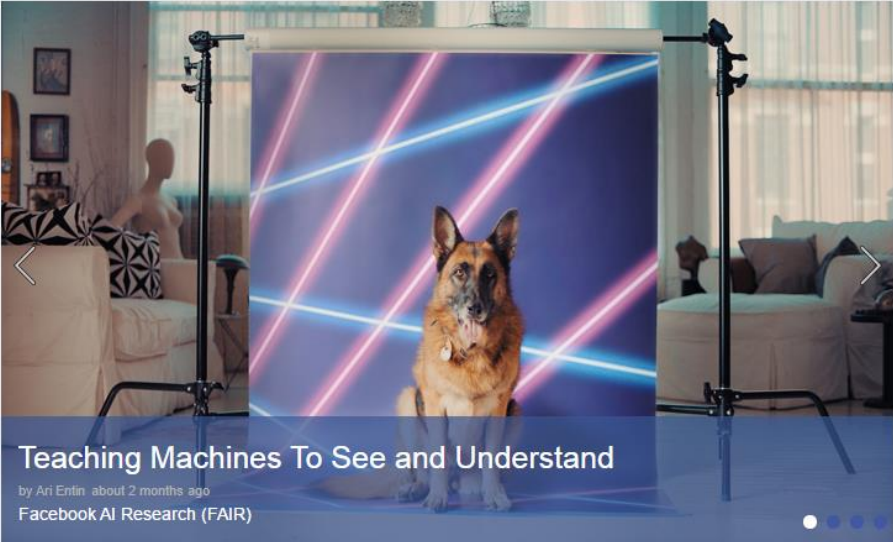
Founded by **Demis Hassabis**, **Shane Legg** and **Mustafa Suleyman** in London, 2010. DeepMind was supported by some of the most iconic tech entrepreneurs and investors of the past decade, prior to being acquired by Google in early 2014 in their largest European acquisition to date.

Industry

- Facebook

Facebook AI Research (FAIR)

Home Publications People Research Downloads Blog



Teaching Machines To See and Understand
by Ari Entin about 2 months ago
Facebook AI Research (FAIR)

Highlights

Teaching Machines To See and Understand
by Ari Entin about 2 months ago
Blog post

Simple bag-of-words baseline for visual question answering
by Bolei Zhou, Yuandong Tian, Sainbayar Sukhbaatar, Arthur Szlam, Rob Fergus about 2 months ago
Publication

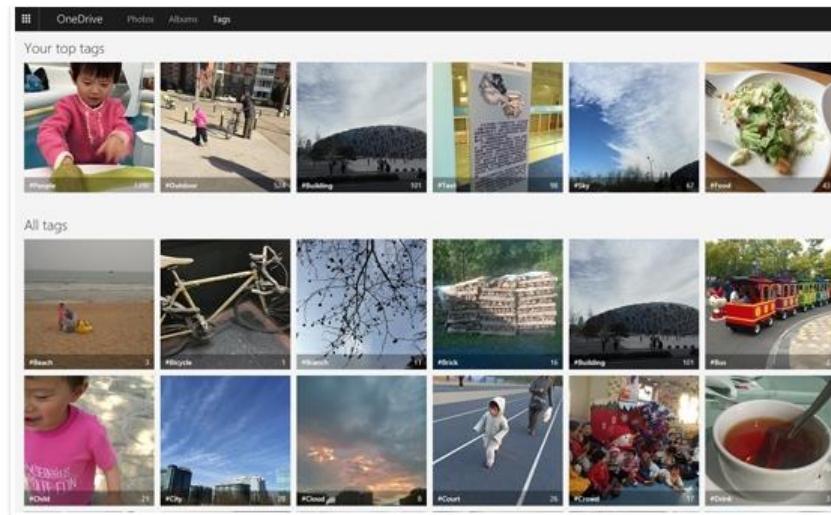
A Roadmap towards Machine Intelligence
by Tomas Mikolov, Armand Joulin, Marco Baroni about 2 months ago
Publication

MazeBase: A Sandbox for Learning from Games
by Sainbayar Sukhbaatar, Arthur Szlam, Gabriel Synnaeve, Soumith Chintala, Rob Fergus about 2

Industry

- Microsoft

Microsoft Researchers' Algorithm Sets ImageNet Challenge Milestone



Industry

- Toyota

SEARCH

The New York Times

TECHNOLOGY

Toyota Invests \$1 Billion in Artificial Intelligence

By JOHN MARKOFF NOV. 6, 2015



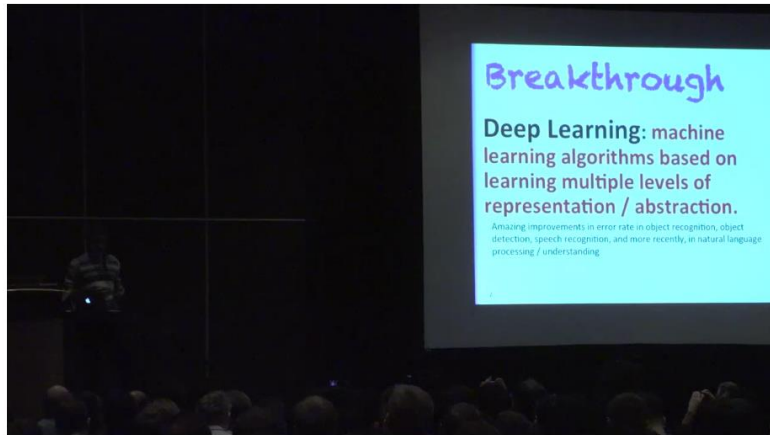
Gill Pratt, a roboticist who will oversee Toyota's new research laboratory in the United States, at a news conference Friday in Tokyo. Yuya Shino/Reuters

Academy

- NIPS 2015: ~4000 attendees, double the number of NIPS 2014



Tutorial: Deep Learning



Academy

- Science special issue
- Nature invited review

REVIEW

Deep learning

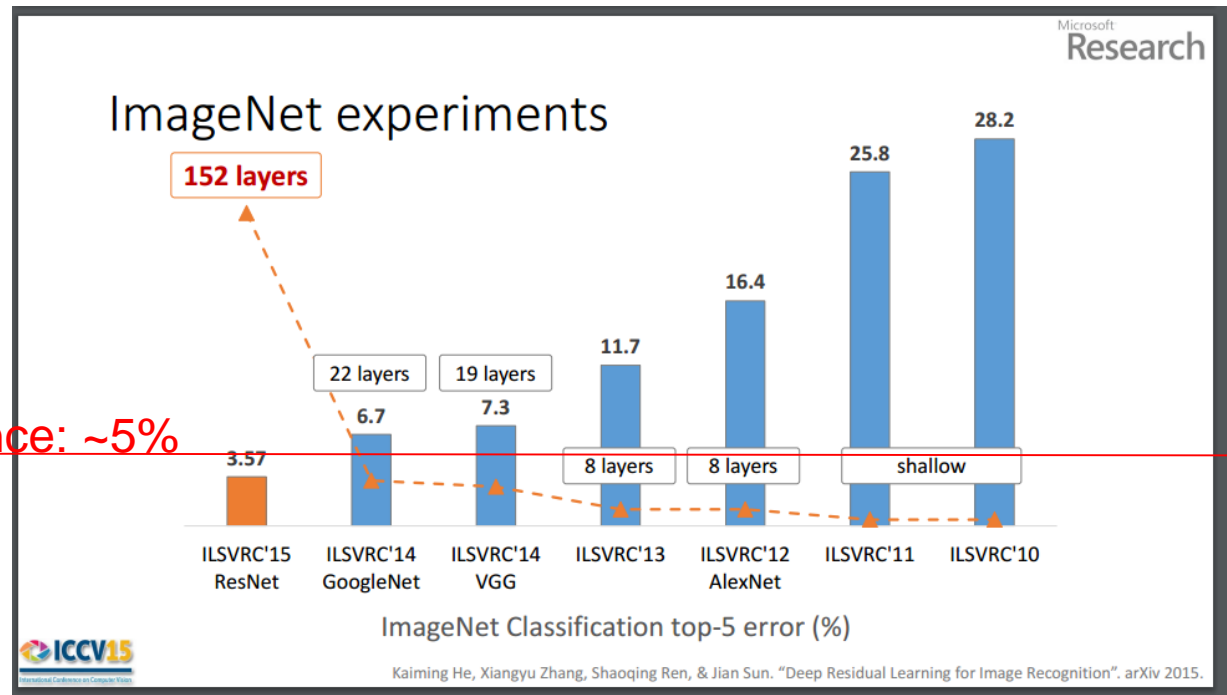
Yann LeCun^{1,2}, Yoshua Bengio³ & Geoffrey Hinton^{4,5}



Image

- Image classification
 - 1000 classes

Human performance: ~5%

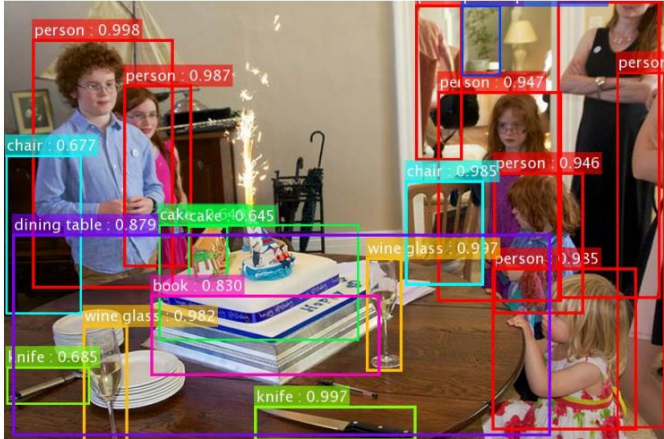


Slides from Kaimin He, MSRA

Image

- Object location

Microsoft Research



person : 0.998
person : 0.987
person : 0.947
person : 0.946
person : 0.935
person : 0.935
chair : 0.677
chair : 0.935
dining table : 0.879
cake : 0.645
wine glass : 0.997
book : 0.830
wine glass : 0.982
knife : 0.685
knife : 0.997

Our results on COCO – too many objects, let's check carefully!

*the original image is from the COCO dataset

ICCV15 International Conference on Computer Vision

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015.
Shaoqing Ren, Kaiming He, Ross Girshick, & Jian Sun. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks". NIPS 2015.

Slides from Kaimin He, MSRA

Image

- Image captioning

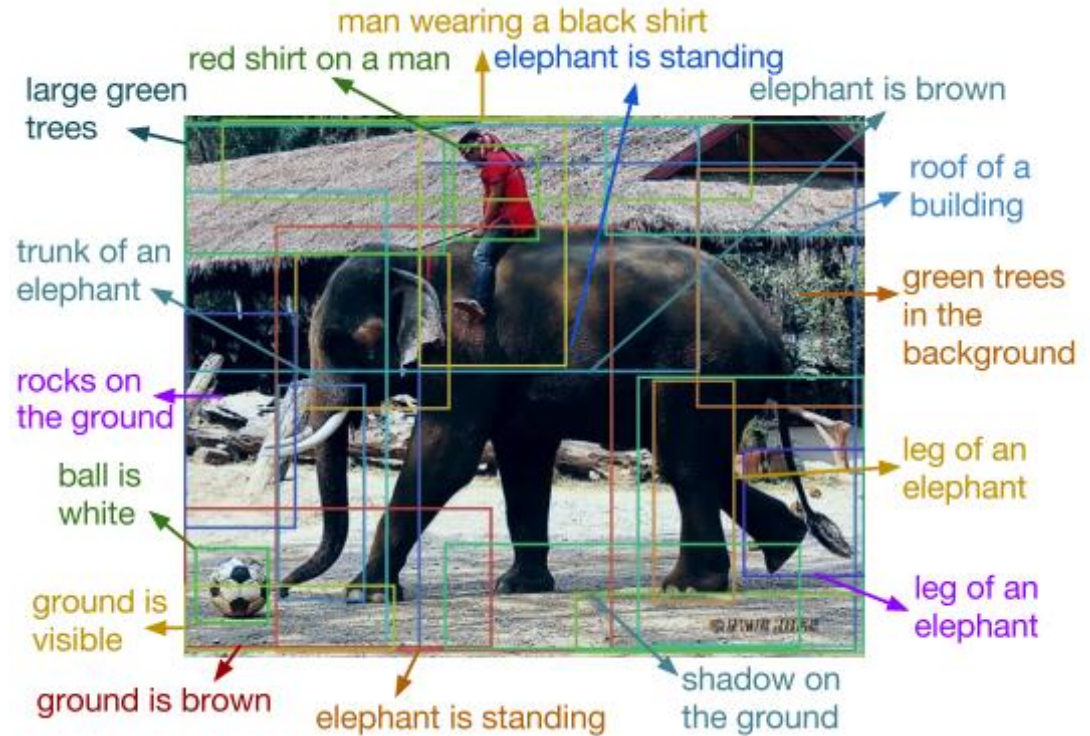


Figure from the paper "DenseCap: Fully Convolutional Localization Networks for Dense Captioning", by Justin Johnson, Andrej Karpathy, Li Fei-Fei

Text

- Question & Answer

I:	Jane went to the hallway.	I:	The answer is far from obvious.
I:	Mary walked to the bathroom.	Q:	In French?
I:	Sandra went to the garden.	A:	La réponse est loin d'être évidente.
I:	Daniel went back to the garden.		
I:	Sandra took the milk there.		
Q:	Where is the milk?		
A:	garden		

Figures from the paper "Ask Me Anything: Dynamic Memory Networks for Natural Language Processing", by Ankit Kumar, Ozan Irsoy, Peter Ondruska, Mohit Iyyer, James Bradbury, Ishaan Gulrajani, Richard Socher

Game



[Google DeepMind's Deep Q-learning playing Atari Breakout](#)

From the paper "Playing Atari with Deep Reinforcement Learning",
by Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou,
Daan Wierstra, Martin Riedmiller

Game



The impact

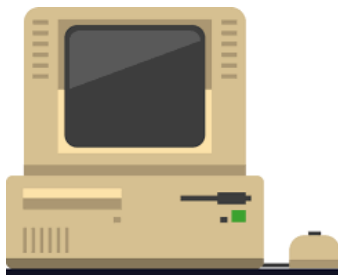
- Revival of Artificial Intelligence
- Next technology revolution?
- A big thing ongoing, should not miss

MACHINE LEARNING BASICS

What is machine learning?

- “A computer program is said to learn from **experience E** with respect to some class of **tasks T** and **performance measure P**, if its performance at tasks in T as measured by P, improves with experience E.”

----- *Machine Learning*, Tom Mitchell, 1997



learning
→



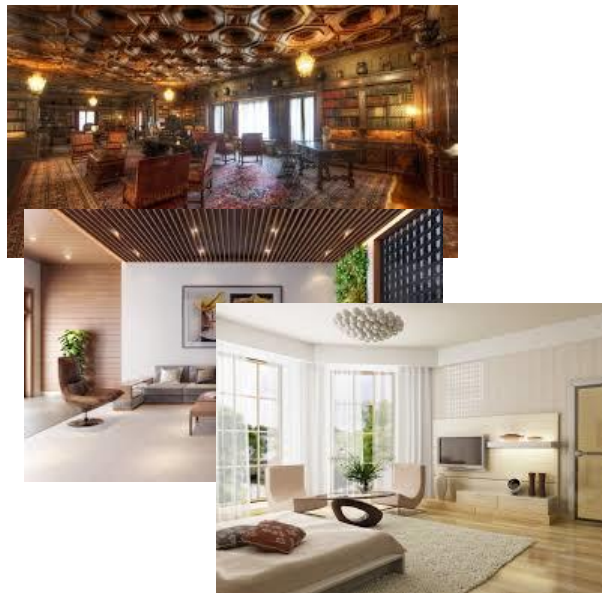
Example 1: image classification



Task: determine if the image is indoor or outdoor

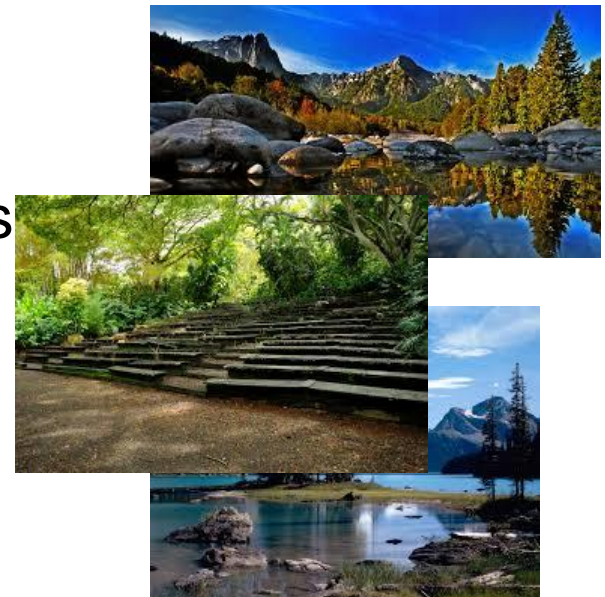
Performance measure: probability of misclassification

Example 1: image classification



Indoor

Experience/Data:
images with labels



outdoor

Example 1: image classification

- A few terminologies
 - Instance
 - Training data: the images given for learning
 - Test data: the images to be classified

Example 1: image classification (multi-class)



ImageNet figure borrowed from vision.stanford.edu

Example 2: clustering images



Task: partition the images into 2 groups
Performance: similarities within groups
Data: a set of images

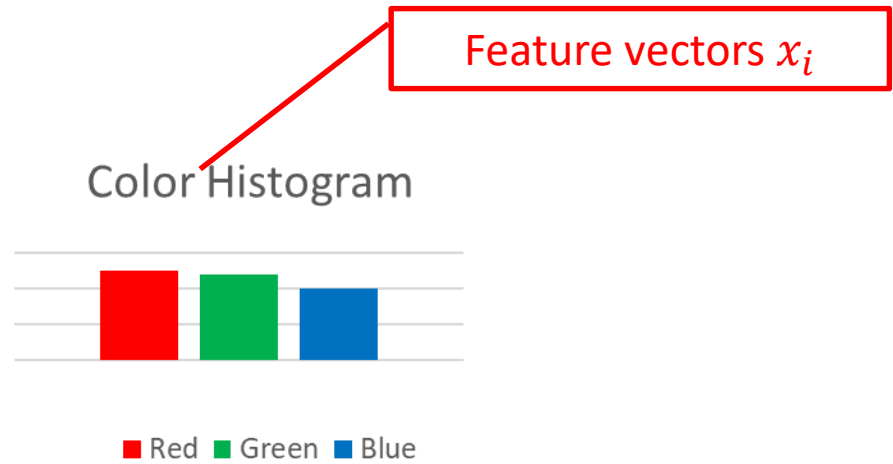
Example 2: clustering images

- A few terminologies
 - Unlabeled data vs labeled data
 - Supervised learning vs unsupervised learning

Feature vectors



Extract
features



Indoor



0

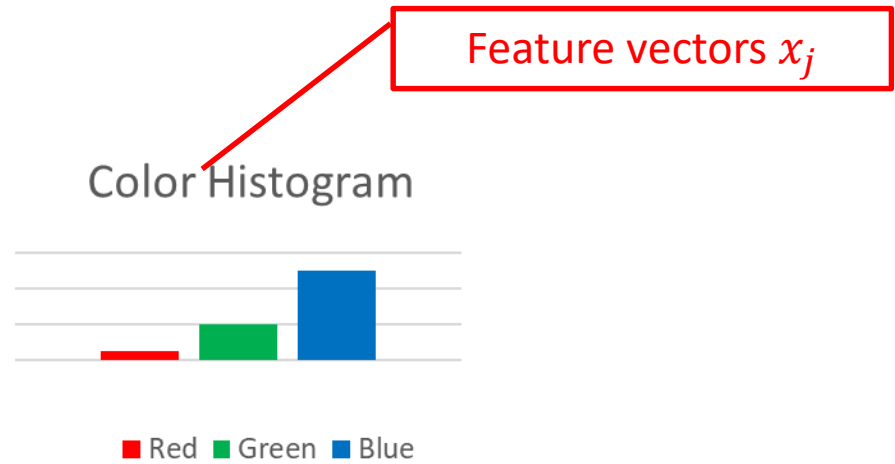
Label y_i

Feature space

Feature vectors



Extract
features



outdoor



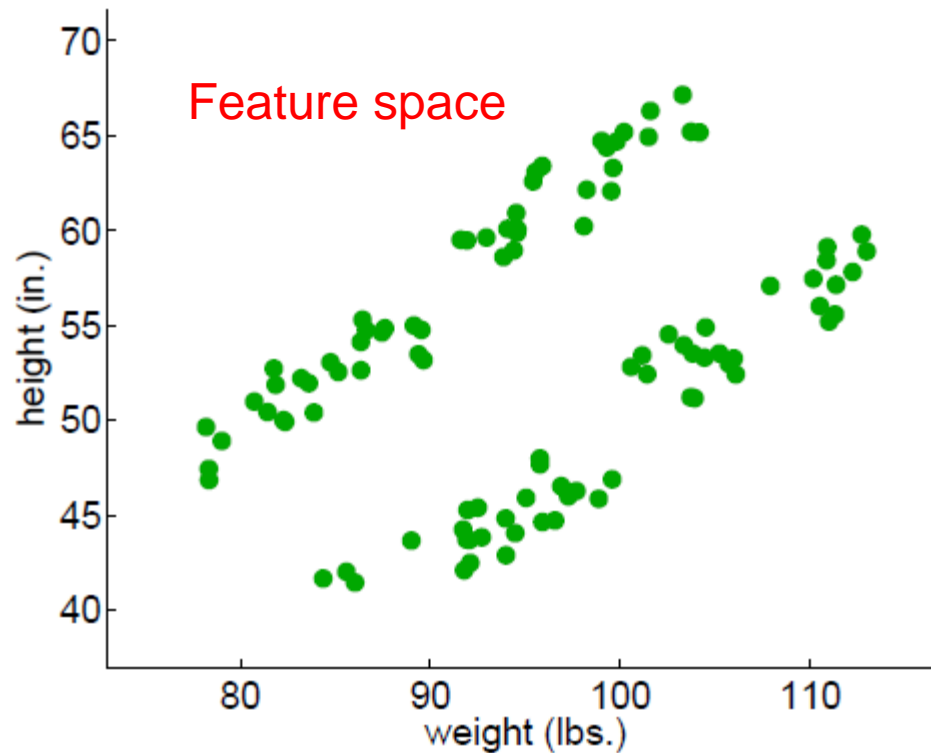
1

Label y_j

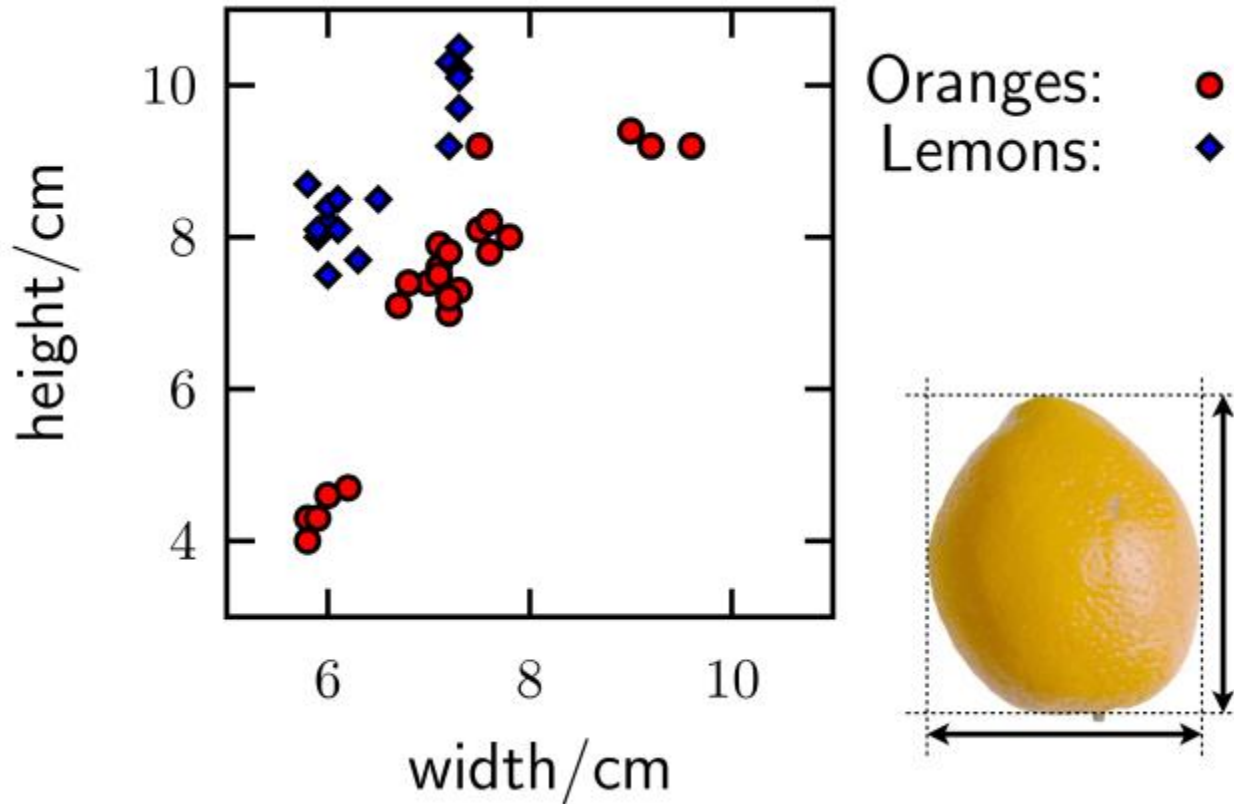
Feature space

Feature Example 2: little green men

- The weight and height of 100 little green men



Feature Example 3: Fruits



- From Iain Murray <http://homepages.inf.ed.ac.uk/imurray2/>

Feature example 4: text

- Text document
 - Vocabulary of size D ($\sim 100,000$)
- “bag of word”: counts of each vocabulary entry
 - To marry my true love \rightarrow (3531:1 13788:1 19676:1)
 - I wish that I find my soulmate this year \rightarrow (3819:1 13448:1 19450:1 20514:1)
- Often remove stopwords: the, of, at, in, ...
- Special “out-of-vocabulary” (OOV) entry catches all unknown words

UNSUPERVISED LEARNING BASICS

Unsupervised learning

in unsupervised learning, we're given a set of instances, **without labels**

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)}$$

goal: discover interesting regularities/structures/patterns that characterize the instances

Common tasks:

- **clustering**, separate the n instances into groups
- **novelty detection**, find instances that are very different from the rest
- **dimensionality reduction**, represent each instance with a lower dimensional feature vector while maintaining key characteristics of the training samples

Anomaly detection

learning
task

given

- training set of instances $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)}$

output

- model h that represents “normal” \mathbf{x}

performance
task

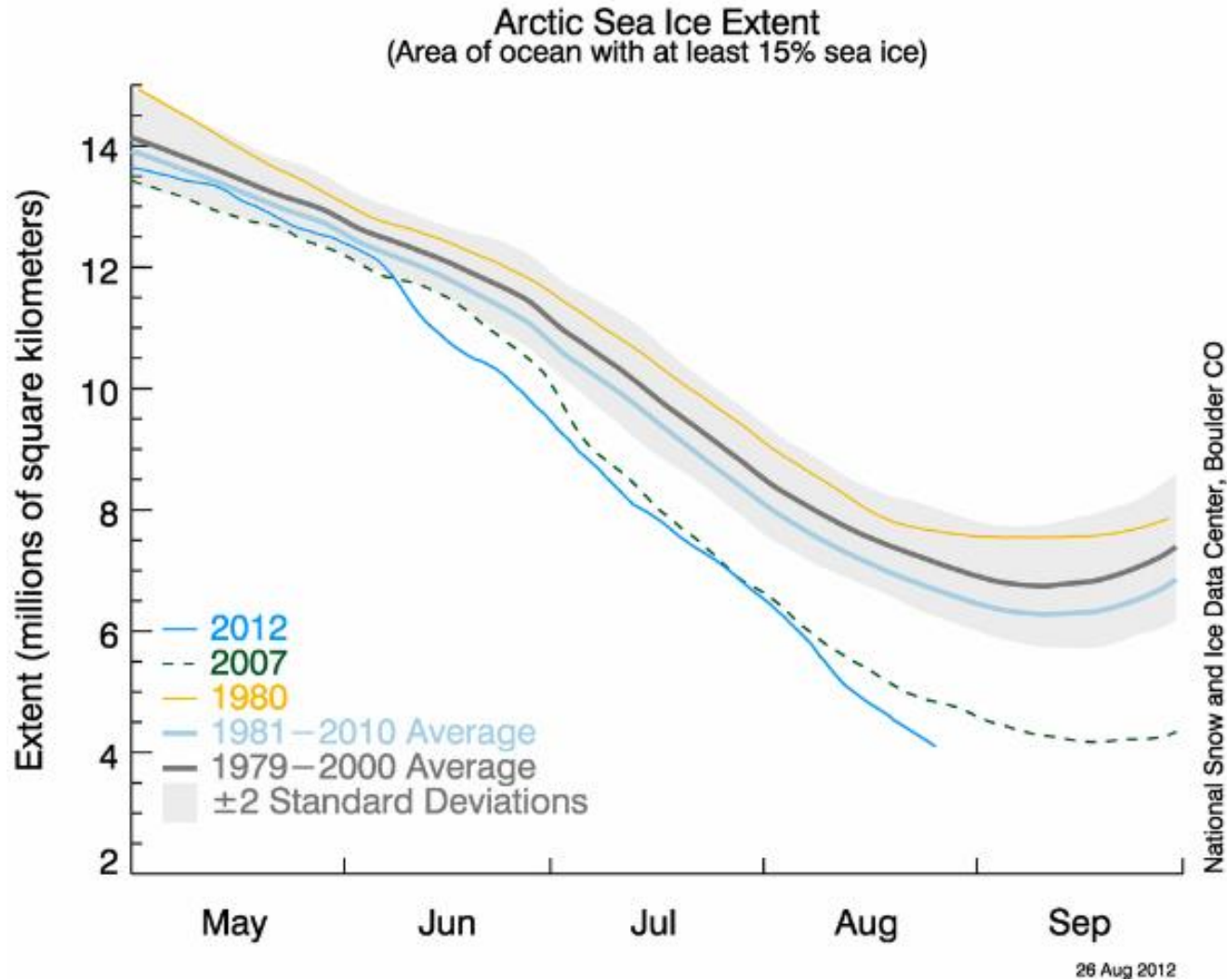
given

- a previously unseen \mathbf{x}

determine

- if \mathbf{x} looks normal or anomalous

Anomaly detection example



Let's say our model is represented by: 1979-2000 average, ± 2 stddev
Does the data for 2012 look anomalous?

Dimensionality reduction

given

- training set of instances $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)}$

output

- model h that represents each \mathbf{x} with a lower-dimension feature vector while still preserving key properties of the data

Dimensionality reduction example



We can represent a face using all of the pixels in a given image

More effective method (for many tasks):
represent each face as a linear
combination of *eigenfaces*



Clustering

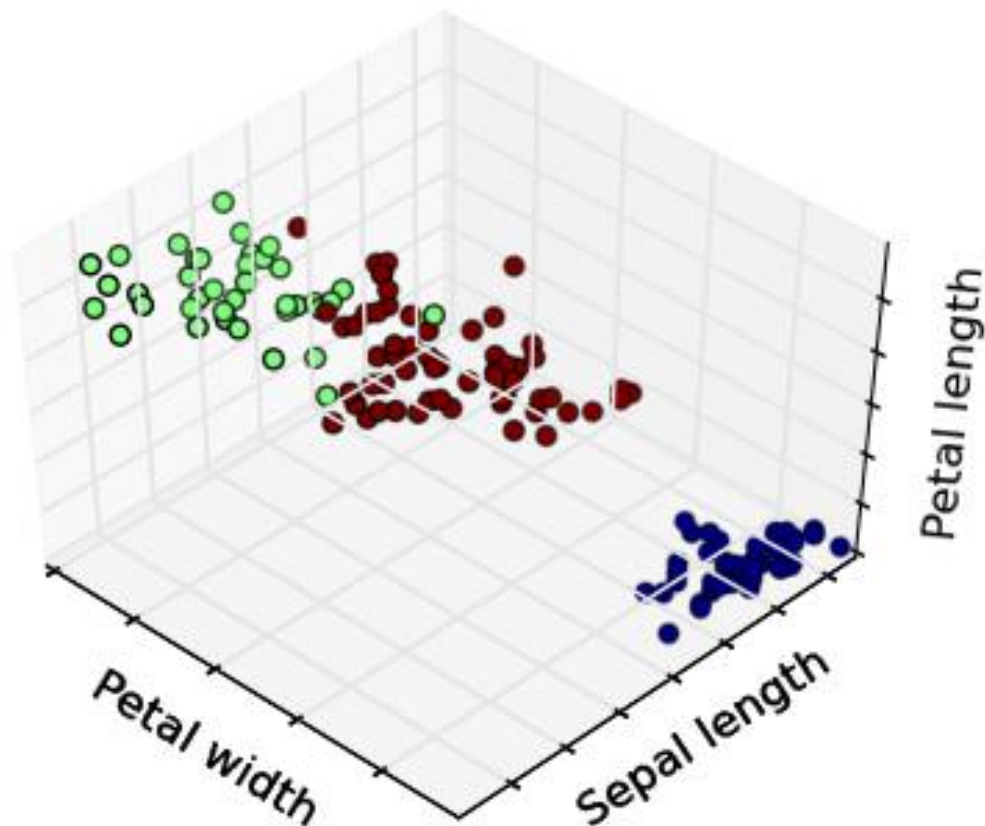
given

- training set of instances $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)}$

output

- model h that divides the training set into clusters such that there is intra-cluster similarity and inter-cluster dissimilarity

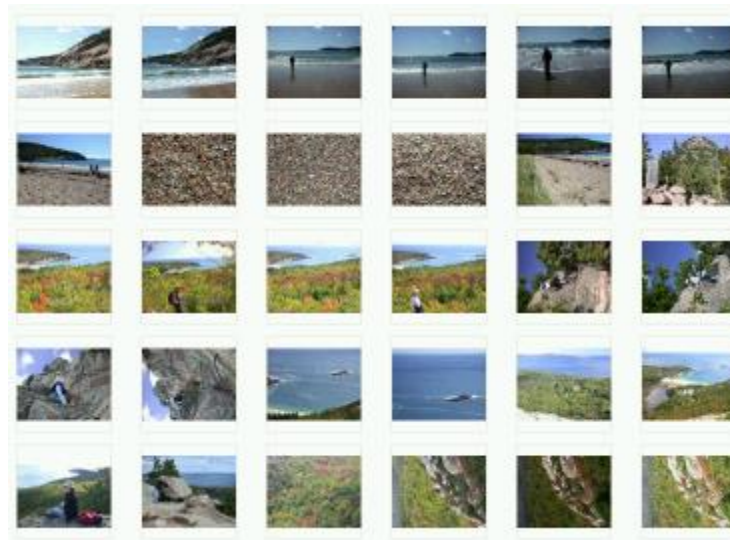
Example 1: Irises



Clustering irises using three different features (the colors represent clusters identified by the algorithm, not y 's provided as input)

Example 2: your digital photo collection

- You probably have >1000 digital photos, 'neatly' stored in various folders...
- After this class you'll be about to organize them better
 - Simplest idea: cluster them using image creation time (EXIF tag)
 - More complicated: extract image features



Two most frequently used methods

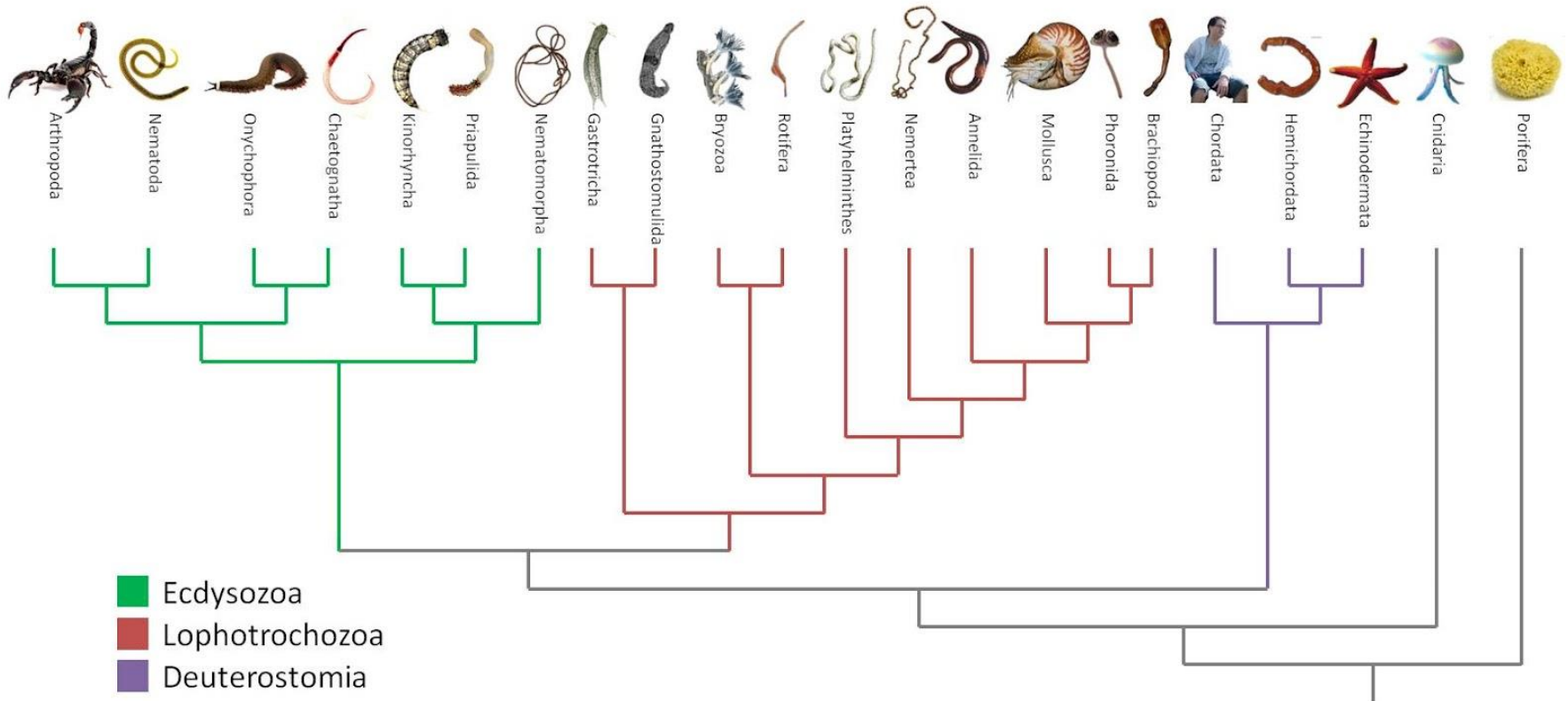
- Many clustering algorithms. We'll look at the two most frequently used ones:
 - Hierarchical clustering
 - Where we build a binary tree over the dataset
 - K-means clustering
 - Where we specify the desired number of clusters, and use an iterative algorithm to find them

HIERARCHICAL CLUSTERING

Hierarchical clustering

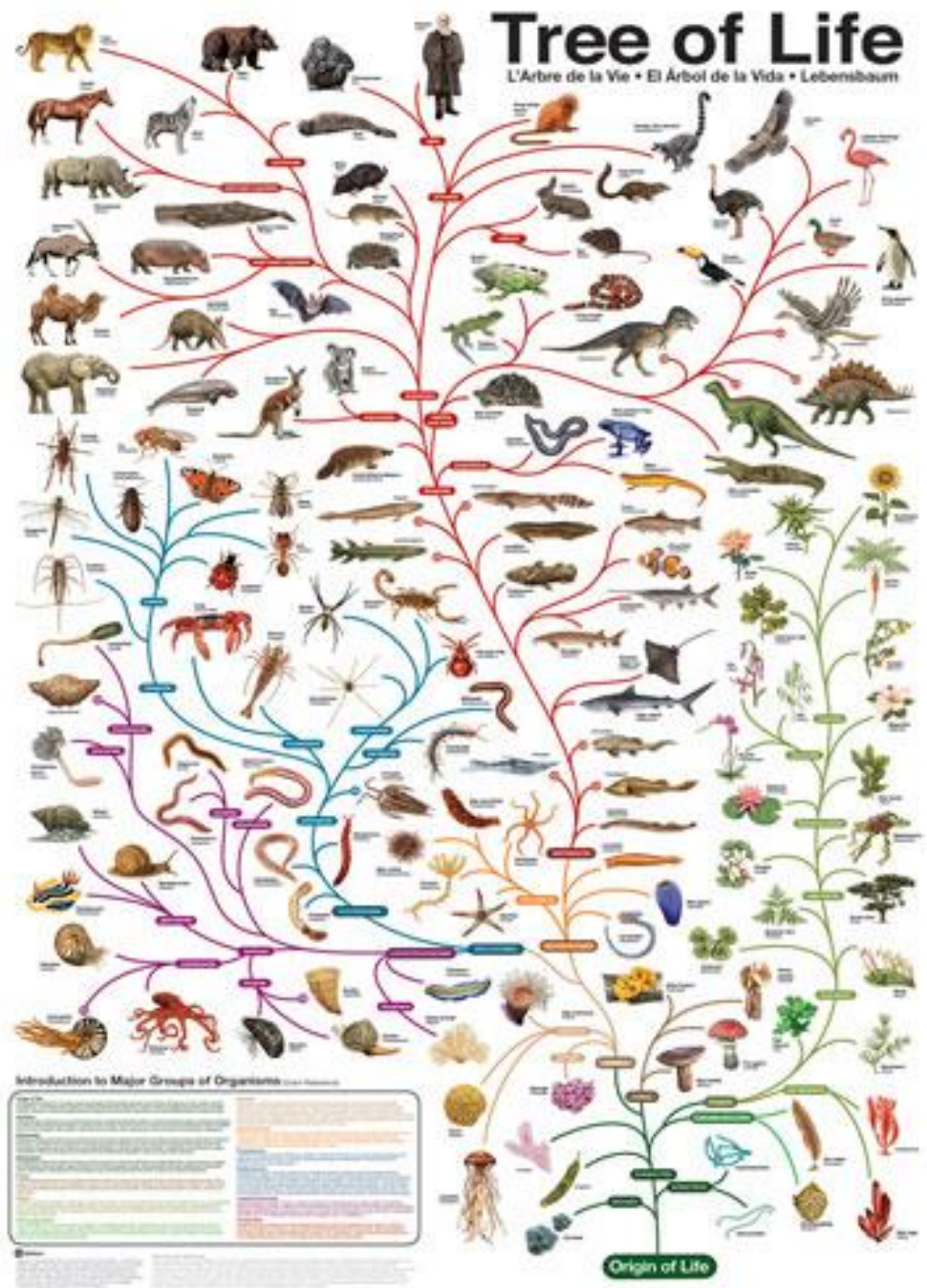
- Very popular clustering algorithm
- Input:
 - A dataset x_1, \dots, x_n , each point is a numerical feature vector
 - Does **NOT** need the number of clusters

Building a hierarchy



Tree of Life

L'Arbre de la Vie • El Árbol de la Vida • Lebensbaum



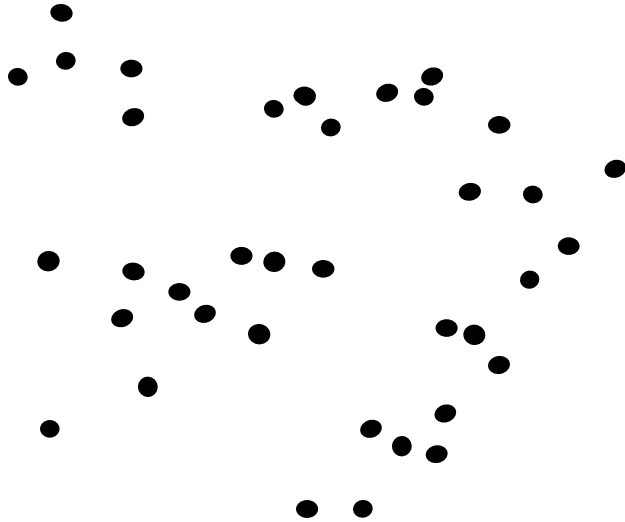
Introduction to Major Groups of Organisms

Group	Characteristics
Bacteria	Prokaryotic, unicellular, diverse metabolism
Archaea	Prokaryotic, often extremophiles
Eukarya	Eukaryotic, diverse forms and functions
- Fungi	Eukaryotic, heterotrophic, cell walls of chitin
- Plantae	Eukaryotic, autotrophic, cell walls of cellulose
- Animalia	Eukaryotic, heterotrophic, no cell walls

© 2004 Pearson Education, Inc. All rights reserved.

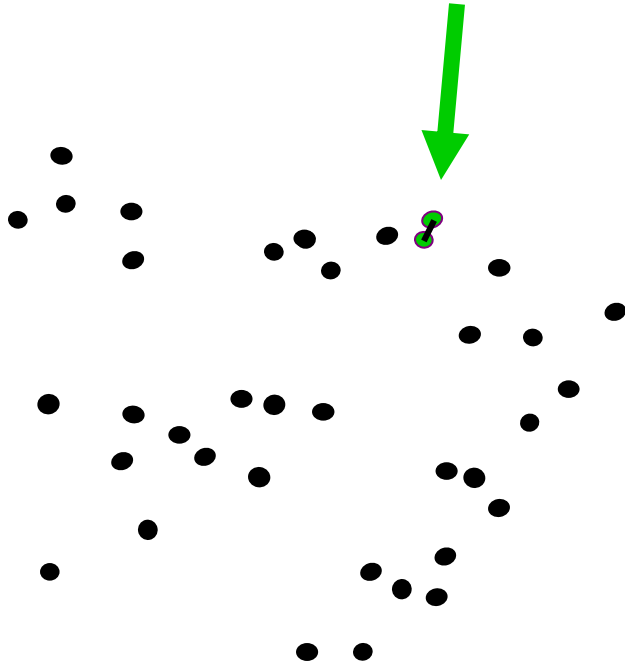
Hierarchical clustering

- Initially every point is in its own cluster



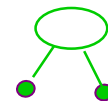
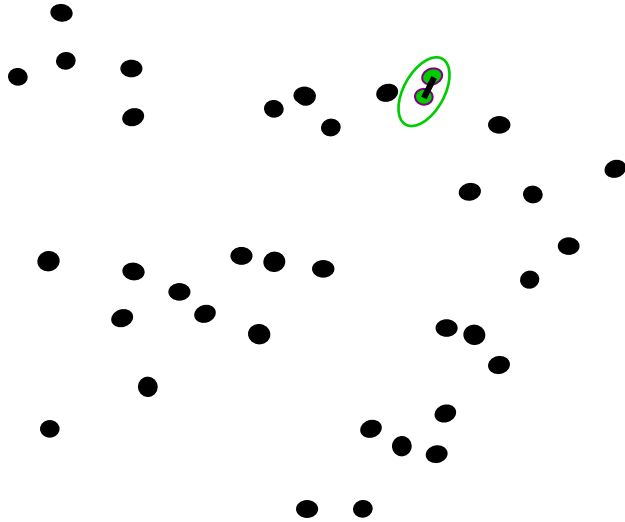
Hierarchical clustering

- Find the pair of clusters that are the closest



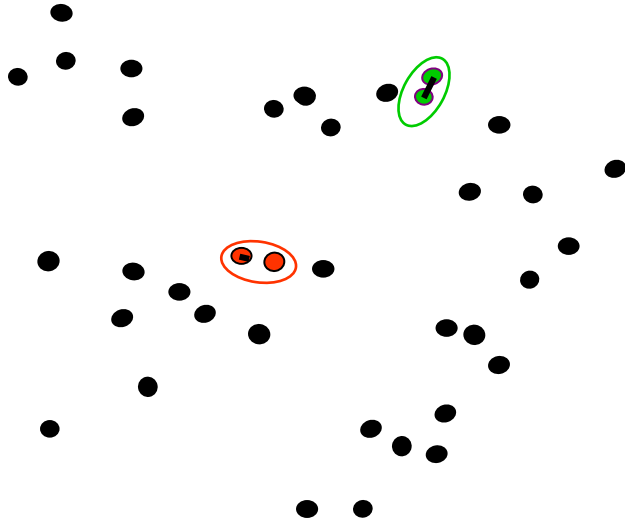
Hierarchical clustering

- Merge the two into a single cluster



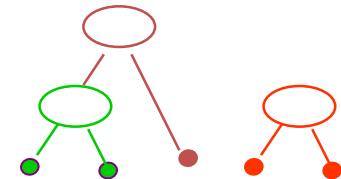
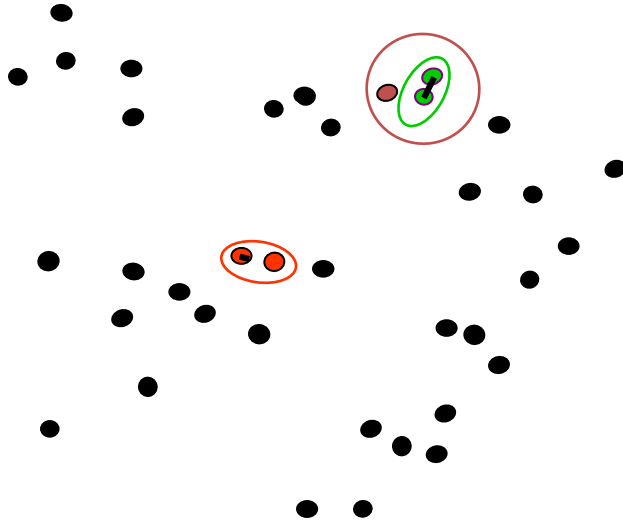
Hierarchical clustering

- Repeat...



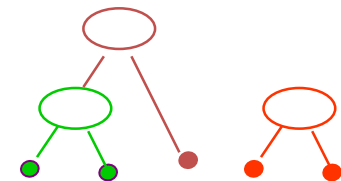
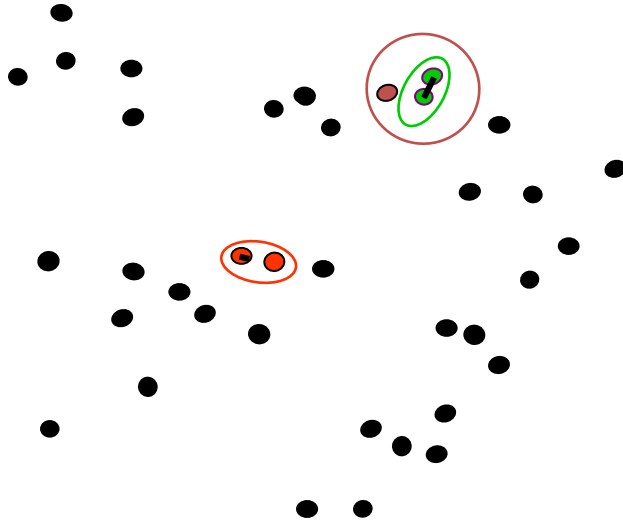
Hierarchical clustering

- Repeat...



Hierarchical clustering

- Repeat...until the whole dataset is one giant cluster
- You get a binary tree (not shown here)



Hierarchical Agglomerative Clustering

Input: a training sample $\{x_i\}_{i=1}^n$; a distance function $d()$.

1. Initially, place each instance in its own cluster (called a singleton cluster).

2. while (number of clusters > 1) do:

3. Find the closest cluster pair A, B , i.e., they minimize $d(A, B)$.

4. Merge A, B to form a new cluster.

Output: a binary tree showing how clusters are gradually merged from singletons to a root cluster, which contains the whole training sample.

- Euclidean (L2) distance

$$d(x_i, x_j) = \|x_i - x_j\| = \sqrt{\sum_{s=1}^d (x_{is} - x_{js})^2}$$

Hierarchical clustering

- How do you measure the closeness between two clusters?

Hierarchical clustering

- How do you measure the closeness between two clusters? At least three ways:
 - **Single-linkage**: the **shortest distance** from any member of one cluster to any member of the other cluster. Formula?
 - **Complete-linkage**: the **greatest distance** from any member of one cluster to any member of the other cluster
 - **Average-linkage**: you guess it!

Hierarchical clustering

- The binary tree you get is often called a **dendrogram**, or **taxonomy**, or a **hierarchy** of data points
- The tree can be cut at various levels to produce different numbers of clusters: if you want k clusters, just cut the $(k - 1)$ longest links
- Sometimes the hierarchy itself is more interesting than the clusters
- However there is not much theoretical justification to it...

K-MEANS CLUSTERING

K-means clustering

- Clustering: What if we want k prototypical examples?

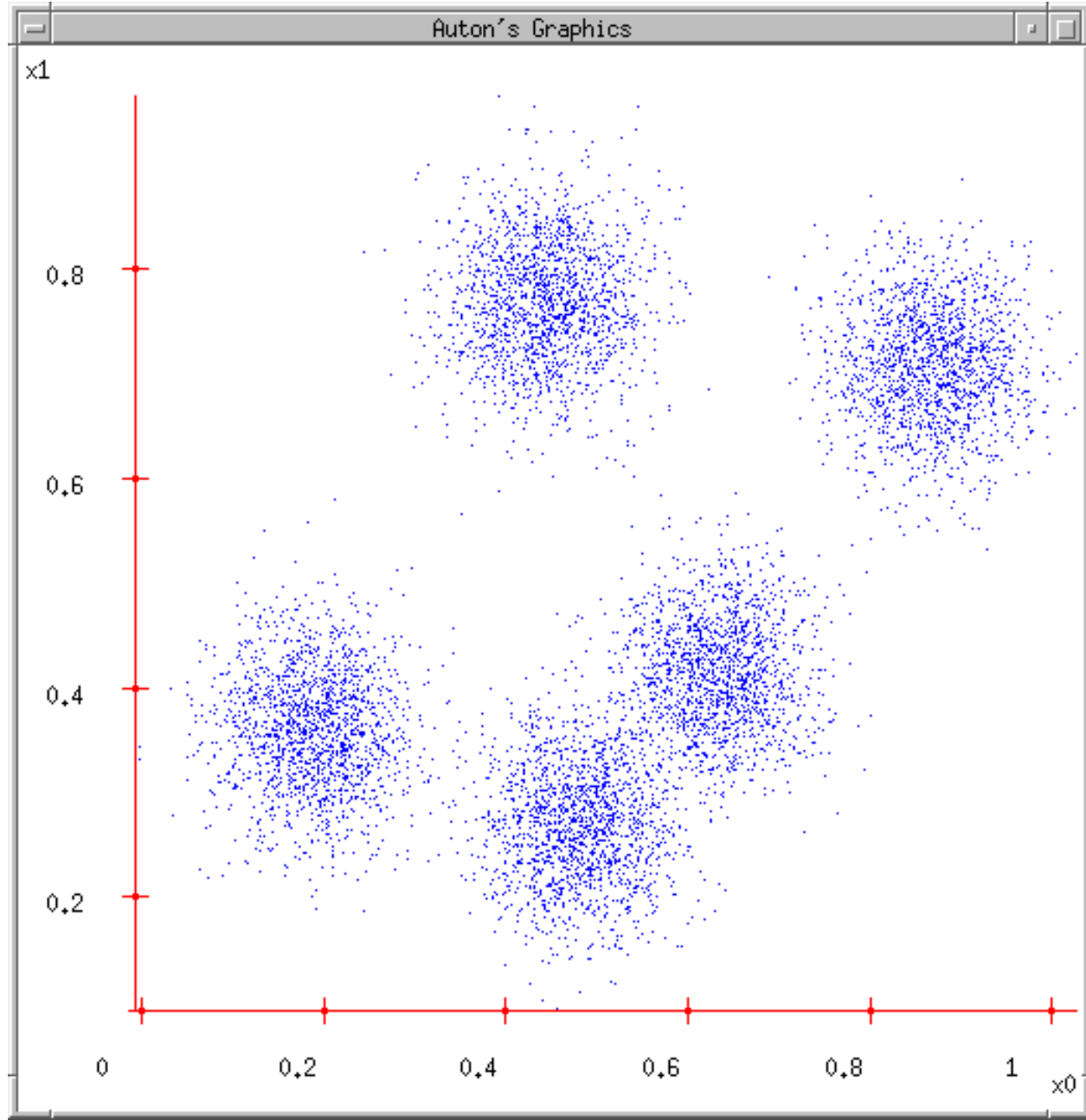


K-means clustering

- Very popular clustering method
- Input:
 - A dataset x_1, \dots, x_n , each point is a numerical feature vector in R^d
 - Assume the number of clusters k is given

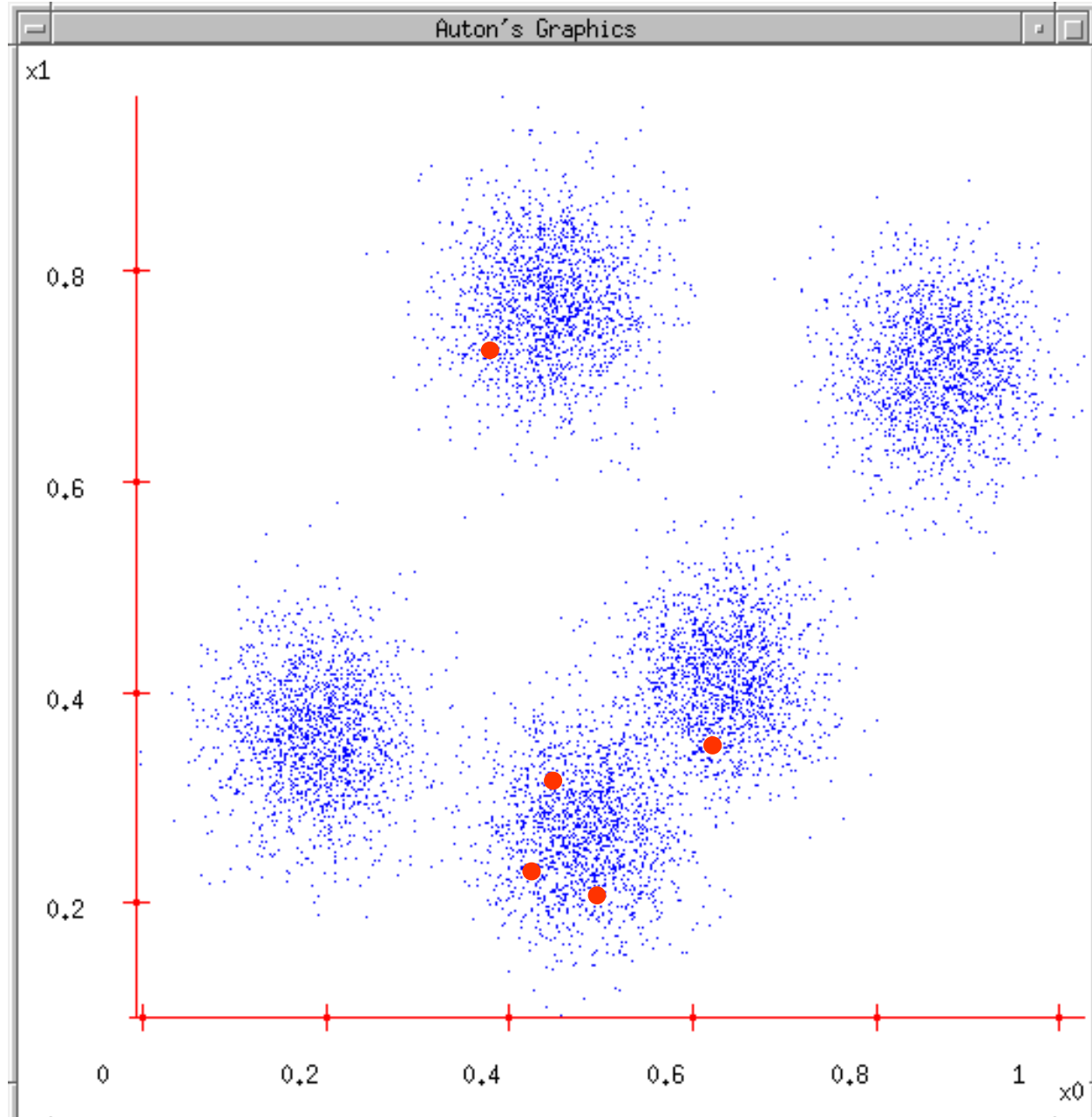
K-means clustering

- Input: dataset, $k = 5$



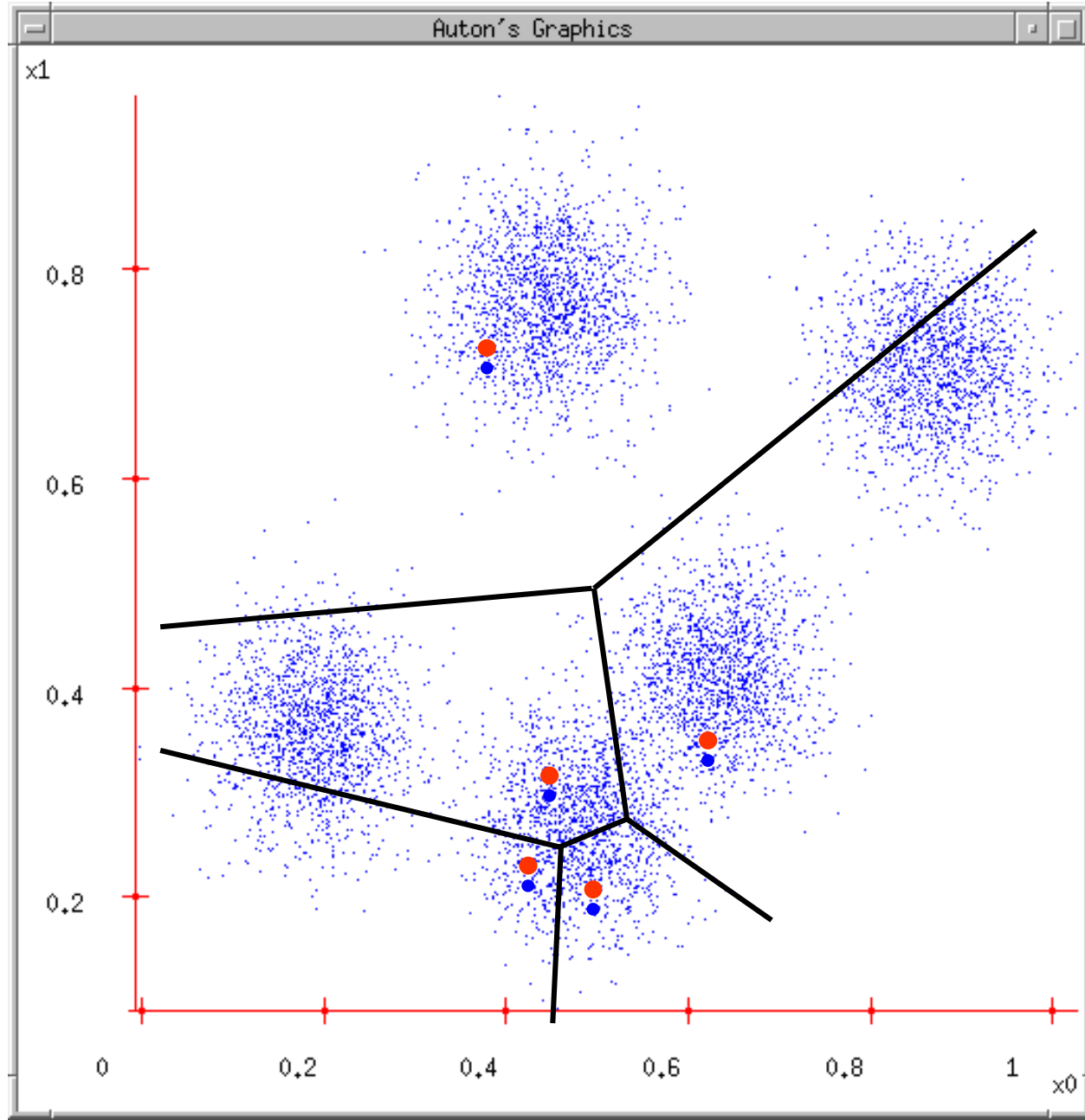
K-means clustering

- Randomly picking 5 positions as initial cluster centers (not necessarily a data point)



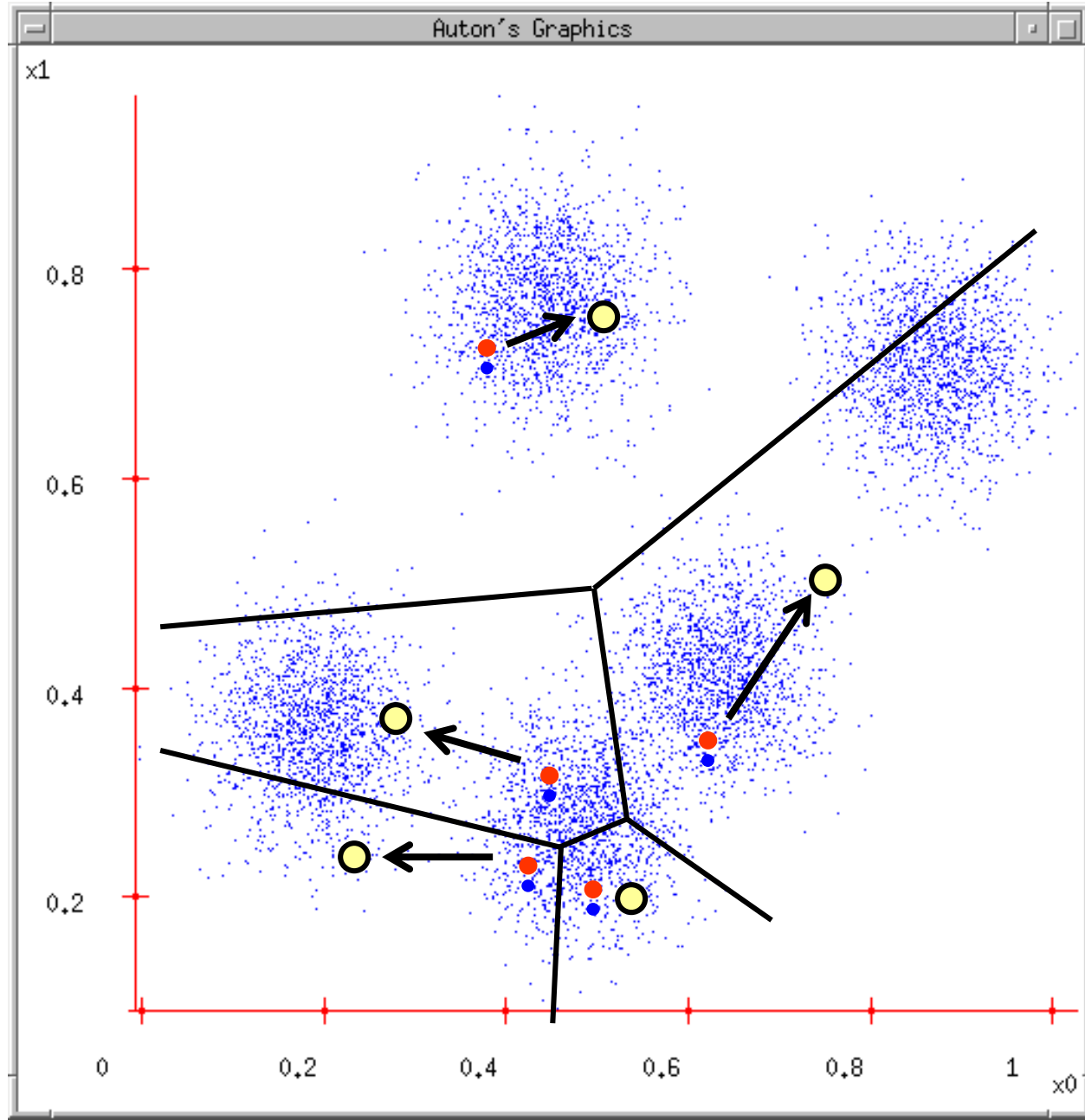
K-means clustering

- Each point finds which cluster center it is closest to. The point is assigned to that cluster.



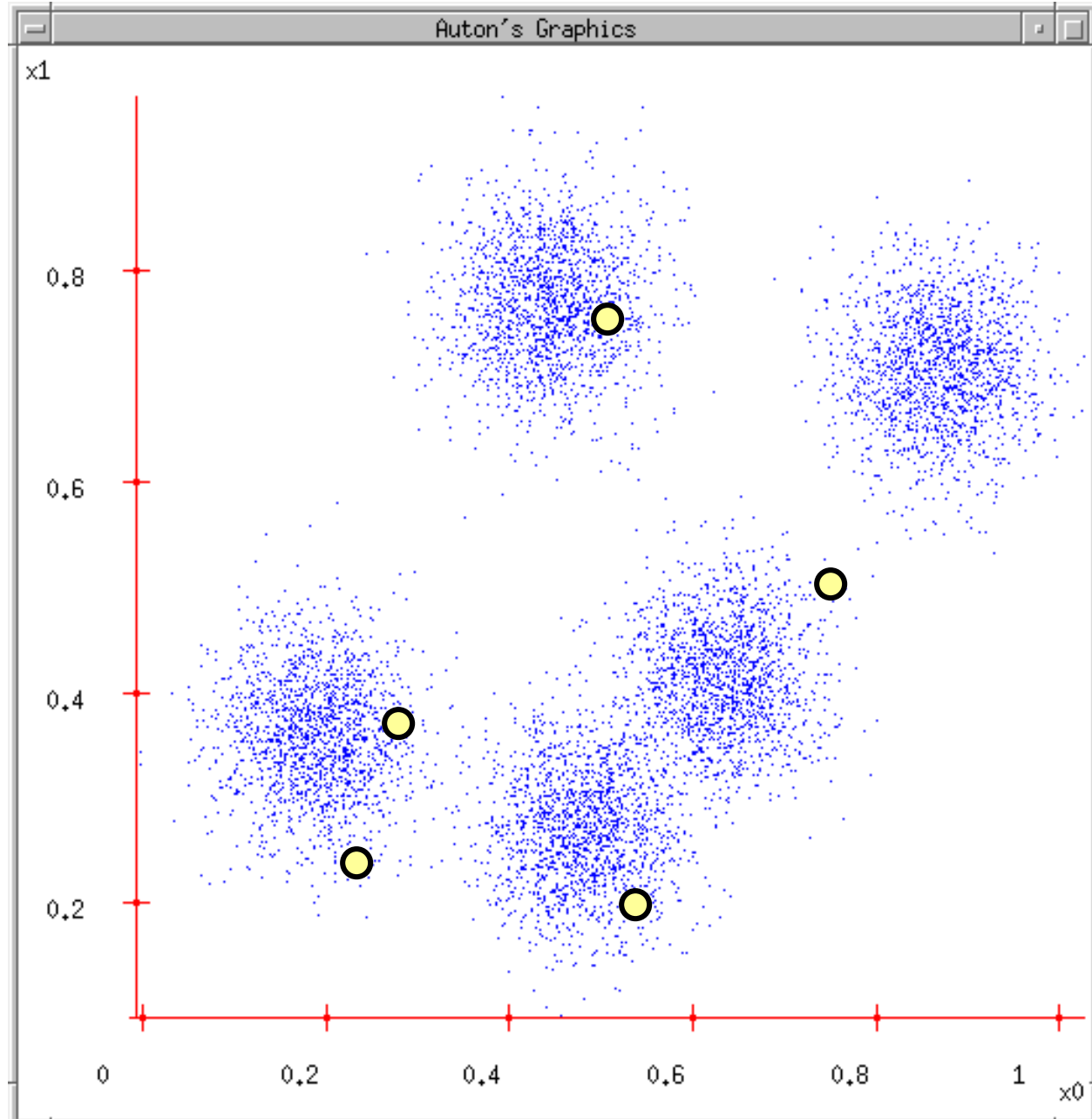
K-means clustering

- Each cluster computes its new centroid, based on which points belong to it



K-means clustering

- Each cluster computes its new centroid, based on which points belong to it
- And repeat until convergence (cluster centers no longer move)...



K-means algorithm

- Input: points x_1, \dots, x_n , number of clusters k
- Select k centers c_1, \dots, c_k
- **Step 1:** for each point x , determine its cluster: find the closest center in Euclidean distance
- **Step 2:** update all cluster centers as the centroids
$$c_i = \sum_{x \text{ in cluster } i} x / \text{SizeOf}(\text{cluster } i)$$
- Repeat step 1, 2 until the centers don't/slightly change

Questions on k-means

- What is k-means trying to optimize?
- Will k-means stop (converge)?
- Will it find a global or local optimum?
- How to pick starting cluster centers?
- How many clusters should we use?

Distortion

- Clustering as summarization: replace a point x with its center $c_{y(x)}$. How far are you off?
- The **distortion** of x is measured by **squared Euclidean distance**:

$$\|x - c_{y(x)}\|^2 = \sum_{i=1}^d [x_i - (c_{y(x)})_i]^2$$

- The distortion of the whole dataset is

$$\sum_x \|x - c_{y(x)}\|^2$$

The optimization objective

- Minimize the distortion of the dataset

$$\min_{\substack{y(x_1), \dots, y(x_n) \\ c_1, \dots, c_k}} \sum_x \|x - c_{y(x)}\|^2$$

Step 1

- Suppose we fix the cluster centers
- Assigning x to its closest cluster center $y(x)$ minimizes the distortion

$$\|x - c_{y(x)}\|^2$$

Step 2

- Suppose we fix the assignment of points. All you can do is to change the cluster centers
- This is a continuous optimization problem!

$$\min_{c_1, \dots, c_k} \sum_x \|x - c_{y(x)}\|^2$$

Step 2

- Suppose we fix the assignment of points. All you can do is to change the cluster centers
- This is a continuous optimization problem!

$$\min_{c_1, \dots, c_k} \sum_x \|x - c_{y(x)}\|^2$$

- Set the gradient to 0 leads to

$$c_i = \frac{\sum_{y(x)=i} x}{n_i}$$

Repeat (step1, step2)

- Both step1 and step2 minimizes the distortion
- Step1 changes the assignments $y(x)$
- Step2 changes the cluster centers c_z

- However there is no guarantee the distortion is minimized over all... need to repeat
- This is hill climbing (coordinate descent)
- Will it stop?

Repeat

- Both step1 and step2 change
- Step1 change
- Step2 change
- However the distortion is minimized
- This is hill climbing
- Will it stop?

There are finite number of points

Finite ways of assigning points to clusters

In step1, an assignment that reduces distortion has to be a new assignment not used before

Step1 will terminate

So will step 2

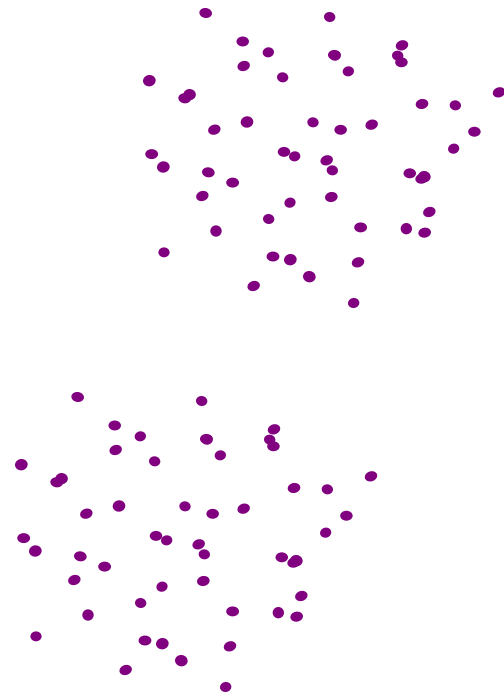
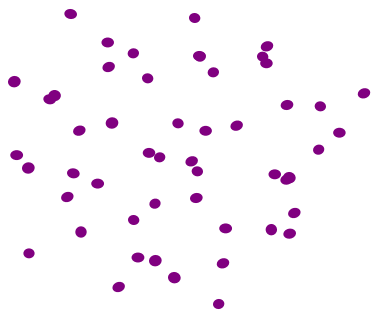
So k-means terminates

Will find global optimum?

- Sadly no guarantee

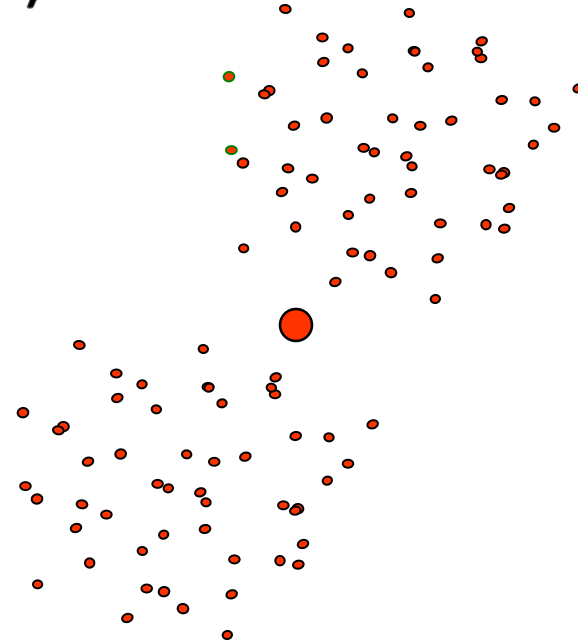
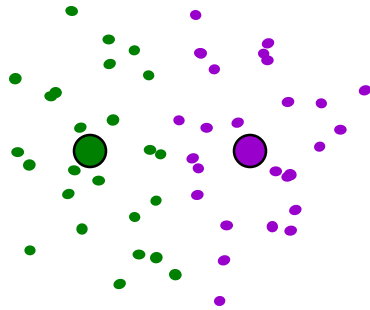
Will find global optimum?

- Sadly no guarantee
- Example (even for $k = 3$)



Will find global optimum?

- Sadly no guarantee
- Example (even for $k = 3$)



Picking starting cluster centers

- Which local optimum k-means goes to is determined solely by the starting cluster centers
 - Be careful how to pick the starting cluster centers. Many ideas. Here's one neat trick:
 1. Pick a random point x_1 from dataset
 2. Find the point x_2 farthest from x_1 in the dataset
 3. Find x_3 farthest from the closer of x_1, x_2
 4. ... pick k points like this, use them as starting centers
 - Run k-means multiple times with different starting cluster centers (hill climbing with random restarts)

Picking the number of clusters

- Difficult problem
- Domain knowledge?
- Otherwise, shall we find k which minimizes distortion?

Picking the number of clusters

- Difficult problem
- Domain knowledge?
- Otherwise, shall we find k which minimizes distortion? $k = n$, distortion = 0
- Need to **regularize**. E.g., the Schwarz criterion

$$\text{distortion} + \lambda(\#param) \log n = \text{distortion} + \lambda dk \log n$$

