

# Machine Learning

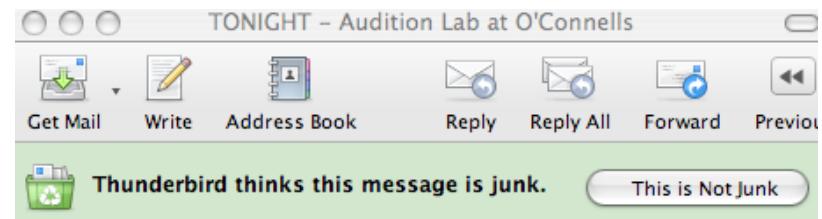
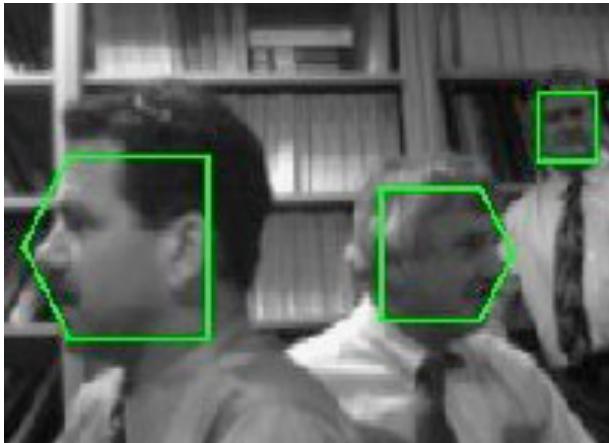
CSCI 4622 Fall 2019

Prof. Claire Monteleoni

# Machine learning is all around us...



# Machine learning is all around us...



DELL COMPUTER 23.01 -1.12	WORLDCOM GRP 14.18 +0.03	PALM INC 12.03 +0.32
------------------------------	-----------------------------	-------------------------

## Gmail

spain



Mail



Move to Inbox



More

59 of about 160



COMPOSE

Group Trips to Spain - [www.madridandbeyond.com/](http://www.madridandbeyond.com/) - Arranging a Group Trip to Spain? Contact Travel Specialist in Madrid

Why this ad?

Inbox  
Chats

Spain excursion



Inbox x

This ad is based on emails from your mailbox. Visit Google's [Ads Preferences](#)

## Gmail

Mail



More

24 of many



COMPOSE

Chegg #1 Textbook Rentals - [www.chegg.com](http://www.chegg.com) - Rent Textbooks & Do The Right Thing Save \$500+ And We Plant a Tree

Why this ad?

Inbox  
Chats

need textbooks back for teaching



Inbox x

This ad is based on emails from your mailbox. Visit Google's [Ads Preferences](#)

## Gmail

Mail



More

9 of many



COMPOSE

Overstock iPads: \$43.20 - [www.NoMoreRack.com/iPad](http://www.NoMoreRack.com/iPad) - Today Only: New 32GB iPads for \$43. 1 Per Customer. Limited Quantities.

Why this ad?

Inbox  
Chats

Your receipt from Apple Store, Georgetown



Inbox x



People (2)

 georgetown@apple.com

to me ▾

12/7/11 georgetown  
Add to circlesSent Mail  
Drafts (39)  
Spam (157)  
[Imap]/Drafts  
[Imap]/Sent



# Gmail

## Mail

 MoreCOMPOSE

68% Off Unlimited Yoga - HomeRun.com/Unlimited\_Yoga\_DC - DC Deals You Won't Find Anywhere! See Today's 3 Unlimited Yoga Deals Why this ad?

InboxChatsSent MailDrafts (39)▼ Important and unread

Woohoo! You've read all the important messages in your inbox.

This ad is based on emails from your mailbox. Visit Google's Ads Preferences Manager to learn more, block specific advertisers, or opt out of personalized ads.



# Gmail

## Mail

 MoreCOMPOSE

Computer Science Degree - cps.Regis.edu/CompScience - Choice Of Specialties, Hands-On & Accelerated Learning. Contact Us! Why this ad?

InboxChatsSent MailDrafts (39)▼ Important and unread

Woohoo! You've read all the important messages in your inbox.

This ad is based on emails from your mailbox. Visit Google's Ads Preferences Manager to learn more, block specific advertisers, or opt out of personalized ads.

# Machine learning

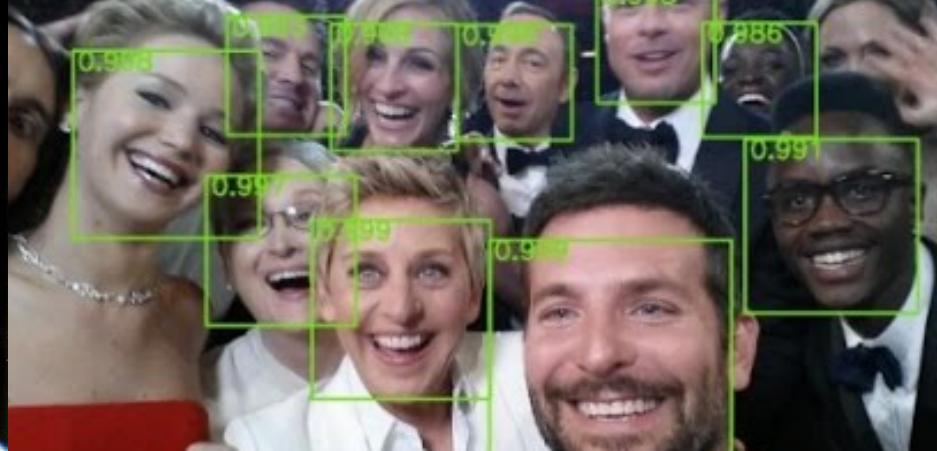
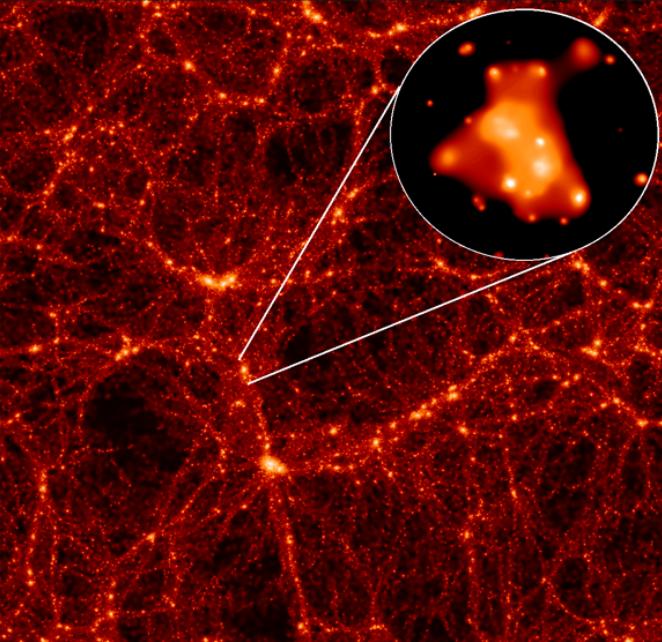
Machine learning is a cutting-edge technology with wide applications, e.g.:

- Web search: Google, Yahoo!, Bing, etc.
- Recommendation systems for books/movies/music, e.g. Amazon, Netflix, Pandora, iTunes
- Personalized internet advertising, e.g. Facebook, Gmail
- Spam filtering
- Autopilot in planes, cruise control in cars
- Finance: Algorithmic trading
- Biology: Bioinformatics
- Chemistry: Computational chemistry
- Big Data – applications in any field



# Take-home messages of course

- Machine learning / AI is not **magic!**
- Yes, ML/AI can be applied to a range of applications!
- **Our goal:** gain conceptual understanding of machine learning algorithms and how to apply them





# We face an **explosion** of data!

Internet transactions, social networks  
DNA sequences  
Financial indicators  
Satellite imagery  
Environmental sensors

....



## Challenges of real data

Massive

many data points

High-dimensional

many features

Noisy, raw

missing data, unlabeled

Sparse

low-dimensional structure

Streaming

time-varying

Sensitive/private

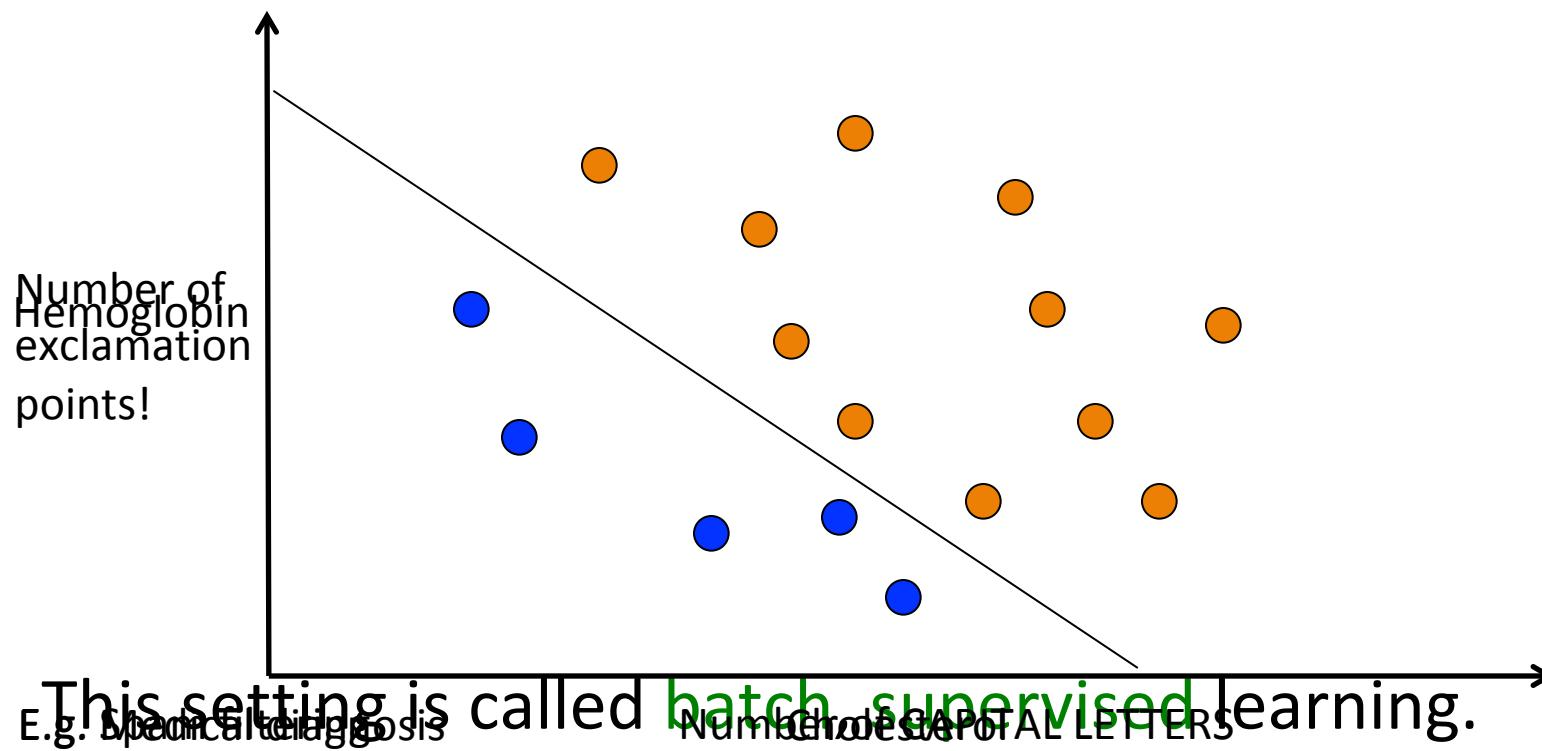
financial, health, social

# Machine Learning 101

Given labeled data points, find a good classification rule.

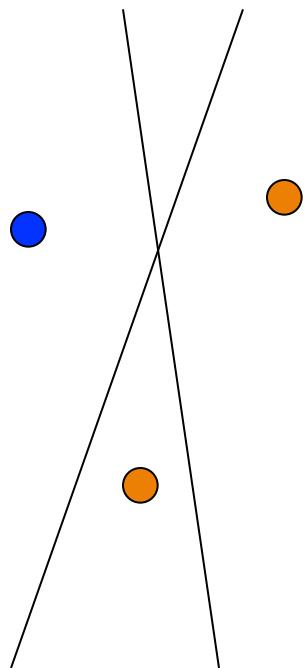
Describes the data

Generalizes well



# Learning from data streams

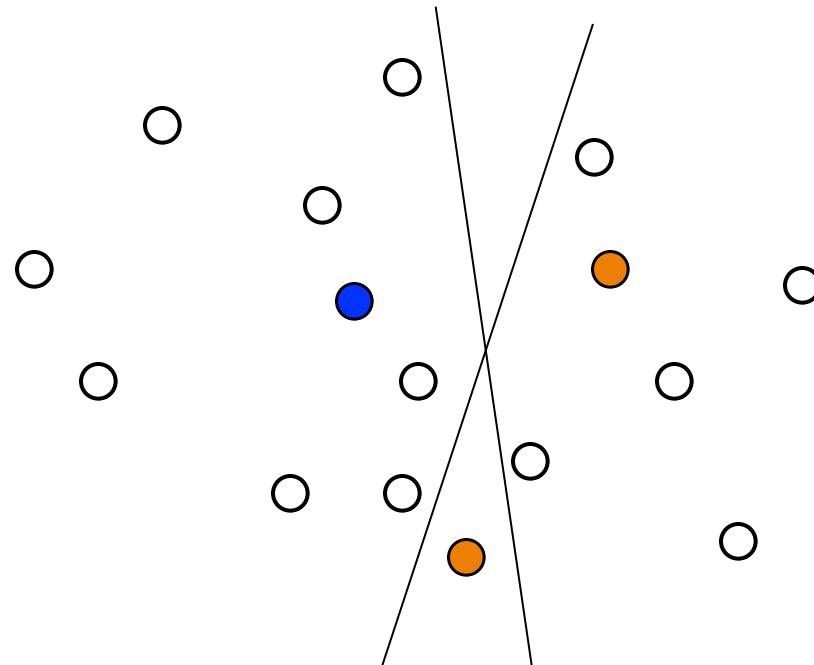
Data arrives in a stream over time.



# Learning from raw data

Semi-supervised learning:

Learning from labeled and unlabeled data

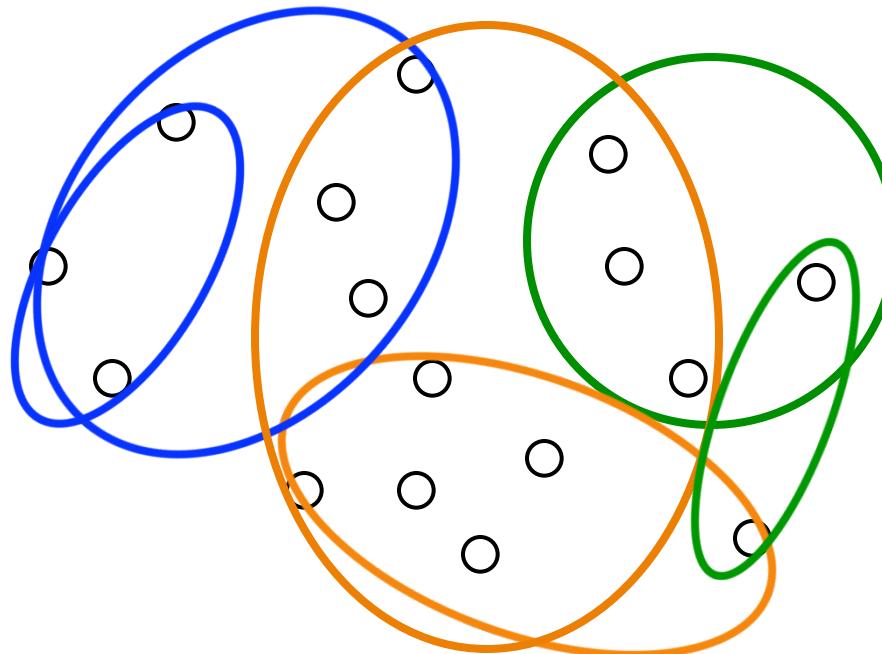


E.g. Active Learning: learner interacts with teacher to query labels.

# Learning from raw data

What can be done without any labels?

Unsupervised learning



E.g. Clustering

# Topics this semester

## Supervised learning

- Introduction to supervised learning (e.g. k-NN, perceptron)
- Generative learning (e.g. HMM)
- Discriminative learning (e.g. SVM, neural networks)
- Regression

## Intro. to learning theory

- Bias-variance tradeoff
- VC dimension

## Unsupervised learning

- Unsupervised deep learning
- Dimensionality reduction / embedding
- Clustering

## Other learning settings

- Online learning
- Active learning

# Course information

Schedule (linked from both Canvas and Piazza)

[sites.google.com/colorado.edu/csci4622fall19/schedule](https://sites.google.com/colorado.edu/csci4622fall19/schedule)

- Please keep up with the readings on the schedule above.
- Please follow **Piazza.com** for important announcements
- Piazza is the preferred mode of communication in this course

TA: Chris Bate

Graduate student staff: Bhalachandra Naik, Aman Satya

My office hours: Tuesdays 2:30-30:30pm, ECES 132

# Piazza

[piazza.com/colorado/fall2019/csci4622/home](https://piazza.com/colorado/fall2019/csci4622/home)  
(linked from Canvas)

PIAZZA CSCI 4622 ▾ Q & A Resources Statistics Manage Class

hw1 hw2 hw3 hw4 hw5 logistics other

Unread Updated Unresolved Following

New Post Search or add a post...

PINNED

Search for Teammates! 8/22/19

LAST WEEK

■ Instr Getting set-up in this course...  
Welcome to Machine Learning! To get set-up for this course you should do the following THIS WEEK: Read the Syllabus

■ Instr Advice for success in this course...  
Welcome to Machine Learning! Advice for success in this course Complete the assigned readings, listed on the Schedule,

■ Private Introduce Piazza to your stu... Thu

Note History: note ★

## Advice for success in this course

Welcome to Machine Learning!

Advice for success in this course

- Complete the assigned readings, listed on the [Schedule](#), before each lecture. Check the schedule for changes/updates.
- Start the Homeworks early. You are encouraged to form study groups of 2-4 students, and work on problems together. Note however, that all your solutions must be your own. (See the Homework Policy in the Syllabus in Canvas).
- Actively participate in the Hands-on sessions, and discuss with your peers.
- Ask and answer questions on Piazza. Helping to answer questions can help you solidify your understanding of the material.

# Course Schedule

[sites.google.com/colorado.edu/csci4622fall19/schedule](https://sites.google.com/colorado.edu/csci4622fall19/schedule)  
(linked from both Canvas and Piazza)

## CSCI 4622: Machine Learning: Fall 2019

Home

Schedule

### Undergrad Machine Learning Fall 2019: Schedule : Sheet1

Week	Date	Lecture	Reading	Assignments (due 5pm)
<b>Week 1</b>	8/26	<u>Introduction, collaborative filtering, nearest neighbor (NN)</u>	[CIML Chapters 1, 3] [AIMA 18.1, 18.2; Mitchell: 8.1, 8.2] [D1]	Intro questionnaire posted.
	8/28	<u>Intro Quiz (on prerequisites: analysis of algorithms, probability, linear algebra), Introduction to linear classification.</u>	[CIML Chapter 4, 6] [Duda, Hart, and Stork: 5.4, 5.5] [D3, J1, J2]	
	8/30	Hands-on 1: Computing Environment Orientation, k-NN [Chris Bate]		Intro questionnaire due.
<b>Week 2</b>	9/4	<u>Perceptron, convergence.</u>	[CIML Chapter 4] [HTF 4.5] [Mitchell: Chapter 3; AIMA: 18.3]	
	9/6			HW 1 posted (on Weeks 1-2)
<b>Week 3</b>	9/9	<u>Decision trees, k-d trees. Overfitting, Bias-variance tradeoff</u>	[CIML Chapters 1, 2][Dasgupta course: Decision Trees] [Description of CART, a canonical DT alg.] [HTF 9.2, Chapters 2, 7]	
	9/11	<u>Bias-variance tradeoff. Model evaluation, cross-validation.</u>	[CIML Chapter 2] [HTF Chapter 7][CIML Chapter 16]	
	9/13	Hands-on 2: Perceptron. Validation.		

# Grading

- 10% Class participation, including on Piazza
- 45% Homeworks
  - 6 assignments: approximately every other week
  - 9% each: your best-scoring 5 out of 6 assignments
- 45% Exams
  - 3 exams (including Final Exam)
  - 25% your top scoring exam, 20% your 2<sup>nd</sup>-top scoring exam
  - Your lowest exam score is dropped
  - Therefore Final Exam is OPTIONAL

# Intro Quiz

- On Wednesday (8/28), class will start with a quiz on the prerequisites:
  - Analysis of Algorithms (CSCI 3104 or equivalent)
  - Probability (CSCI 3022 or equivalent)
  - Linear Algebra (CSCI 2820 or equivalent)
- This will **not** count towards your grade. It is meant to help you decide if you have the preparation.
- If you have taken the prerequisites, the quiz will be straightforward.
- The actual content in this course will be more advanced.

# Other set-up tasks

- Set-up a free account on [www.Piazza.com](http://www.Piazza.com) and enroll in the course. All announcements will be there!
- [Intro Questionnaire](#): please fill this out within the week of the semester (linked from Piazza and Canvas).
- Make sure you have access to Python and Jupyter notebooks, and brush up on basic programming.
- Look at the list of recommended textbooks and make sure you have access. Some are free online, others in the library.
- Do the readings listed in the schedule, before each lecture.
- Form a study group (2-4 students).

# Electronic device policy

- No cell phone / smartphone use (please turn it off before entering the classroom)
- Tablets may be used when kept flat on the desk
- Laptops or tablets with vertical or slanted screens are restricted to the back 3 rows of the lecture hall, to avoid distracting other students. **NOTE: This was a specific student request to CU CEAS faculty.**
- When in doubt, if your screen could be viewed by other students, and cause a distraction, then please sit towards the back of the room.
- Exceptions:
  - **Exams:** no electronic devices allowed
  - **Hands-on days:** please bring your laptop!

# Homework collaboration policy

Collaboration is allowed, and encouraged, to the following extent:

- you are welcome to discuss problems (both written and programming) with each other and to take your own notes during these discussions.
- However, you must write the solutions and implement code on your own;
- **Copying is not permitted.**
- You must write, on the assignment, the names of students you discussed each problem with, and any external sources you used in solving the problem.

# Lateness policy

- Assignments are due Friday at 5pm.
- Late submissions received from 5:01PM Friday until 11:59PM Sunday night will be accepted, with a deduction of 25% credit.
- After that, late submissions will not be accepted. Recall: your lowest HW score will be dropped.

# Advice for success in this course

- Complete the assigned readings, listed on the [Schedule](#), before each lecture. Check the schedule often for changes/updates.
- Start the homeworks early. You are encouraged to form study groups of 2-4 students, and discuss the homework problems together. Note however, that all your solutions must be your own. (See the Homework Collaboration Policy).
- Ask and answer questions on [Piazza](#). Helping to answer questions can help you solidify your own understanding.

# Collaborative Filtering and Nearest Neighbor



## Jade Yoga Elite Yoga Mat

by [Jade Yoga](#)

(1 customer review) | (0)



### Customers Who Bought This Item Also Bought

Page 1 of 4



Aurorae Organic Yoga Mat Wash-Cleans, Restores, Refreshes, Refreshes... by Aurorae

(27)

\$11.95



Jade Harmony Professional 3 / 16-Inch Yoga Mat by Jade Yoga

(162)

\$62.95 - \$105.00



Jade Fusion 74-inch Yoga Mat by Jade Yoga

(36)

\$116.78 - \$169.95

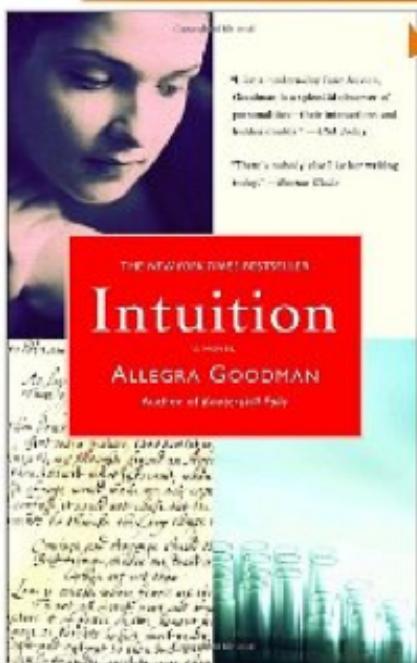


Yoga Mat Bag Extra Large Easy Open Zipper 100% C... by Bean Products

(44)

\$18.95 - \$21.95

Click to **LOOK INSIDE!**



## Intuition [Paperback]

[Allegra Goodman](#)  (Author)

 [\(69 customer reviews\)](#) |  [Like](#) (4)

This intimate portrait of life in a research institute in Cambridge, Massachusetts, revolves around a scientific mystery....

### Customers Who Bought This Item Also Bought

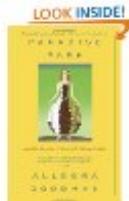
Page 1 of 25



[Kaaterskill Falls](#) by  
Allegra Goodman

 (62)

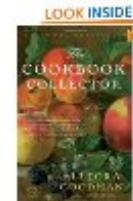
\$10.40



[Paradise Park](#) by Allegra  
Goodman

 (36)

\$11.12



[The Cookbook  
Collector: A Novel](#) by  
Allegra Goodman

 (135)

\$10.20



[Statistics for the Life  
Sciences \(4th Edition\)](#) by  
Myra L. Samuels

 (20)

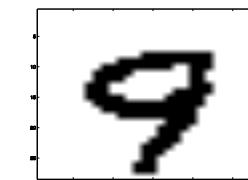
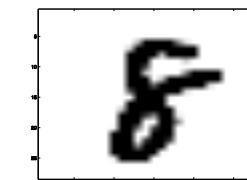
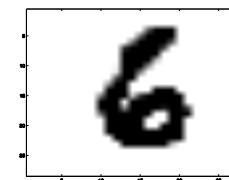
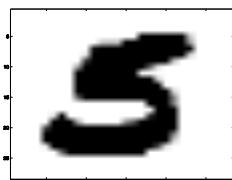
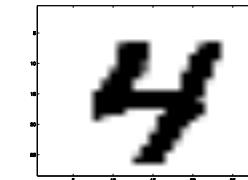
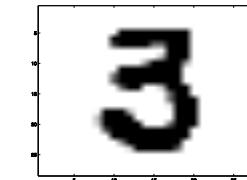
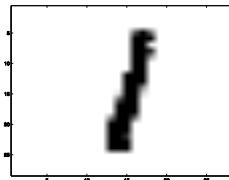
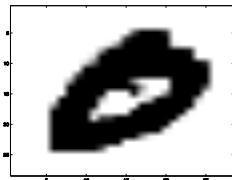
\$113.99

# Nearest Neighbor Classification

[With much credit to S. Dasgupta]

# A first example

Database of 20,000 images of handwritten digits, each labeled by a human



[28 x 28 greyscale; pixel values 0-255; labels 0-9]

Use these to learn a classifier which will label digit-images automatically...



# Data representation

Image: size  $28 \times 28 = 784$  pixels, each pixel is grayscale, in 0-255.

- Unfold each image ( $28 \times 28$  pixels) into a vector of length 784.



- Now we have vectors in 784-dimensional Euclidean space:

$$\text{Data space } X = \{0, 1, \dots, 255\}^{784} \subset \mathbb{R}^{784}$$

- Each image comes with a label, which is a digit.

Labels are in the set  $Y = \{0, 1, \dots, 9\}$

# Training data set

1 4 1 0 1 1 9 1 5 4 8 5 7 2 6 8 0 3 2 2 6 4 1 4 1  
8 6 6 3 5 9 7 2 0 2 9 9 2 9 9 7 2 2 5 1 0 0 4 6 7  
0 1 3 0 8 4 1 1 1 5 9 1 0 1 0 6 1 5 4 0 6 1 0 3 6  
3 1 1 0 6 4 1 1 1 0 3 0 4 7 5 2 6 2 0 9 9 7 9 9  
6 6 8 9 1 2 0 8 6 7 0 8 5 5 7 1 3 1 4 2 7 9 5 5 4  
6 0 1 0 1 8 7 3 0 1 8 7 1 1 2 9 9 3 0 8 9 9 7 0 9  
8 4 0 1 0 9 7 0 7 5 9 7 3 3 1 9 7 2 0 1 5 5 1 9 0  
5 5 1 0 7 5 5 1 8 2 5 5 1 8 2 8 1 4 3 5 8 0 9 0 9  
4 3 1 7 8 7 5 4 1 6 5 5 4 6 0 3 5 4 6 0 3 5 4 6 0  
5 5 1 8 2 5 5 1 0 8 5 0 3 0 4 7 5 2 0 4 3 9 4 0 1

- Each image comes with a label, which is a digit.

Labels are in the set  $Y = \{0, 1, \dots, 9\}$

# The learning problem

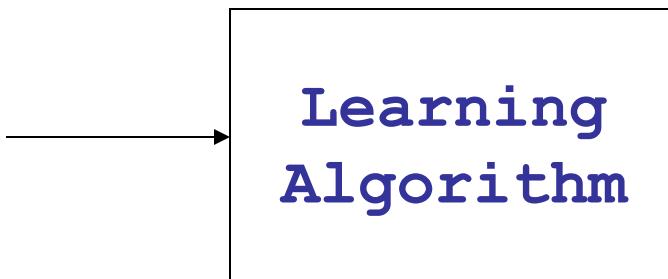
Input space  $X = \{0, 1, \dots, 255\}^{784}$

*f is function for training set*

Output space  $Y = \{0, 1, \dots, 9\}$

*scalar*

Training set  
 $(x_1, y_1), \dots, (x_n, y_n)$   
 $n = 20000$



*we can't learn from single data set*

*what is learn in  
machine learning to produce  
output digit (0 → 9)*

To measure how good  $f$  is: use a test set.

We have a test set with 100 instances of each digit.

# A possible strategy

Input space  $X = \{0, 1, \dots, 255\}^{784}$

Output space  $Y = \{0, 1, \dots, 9\}$

Treat each image as a point in 784-dimensional Euclidean space

To classify a new (unlabeled) image:

find its nearest neighbor in the database (training set)  
and return that label

**f = search engine + entire training set**

# Euclidean distance in high dimension

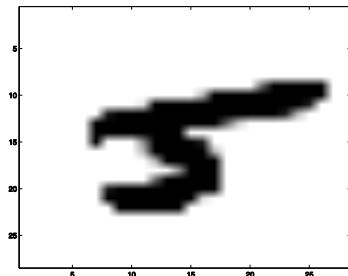
Euclidean distance between 784-dimensional vectors  $x, z$  is

$$\|x - z\| = \sqrt{\sum_{i=1}^{784} (x_i - z_i)^2}$$

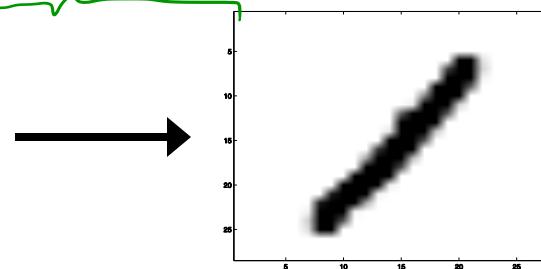
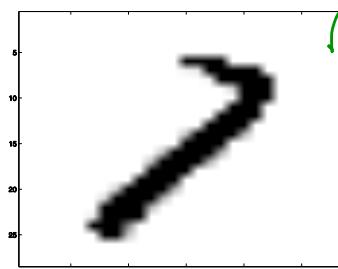
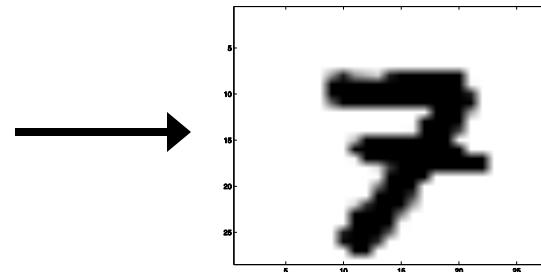
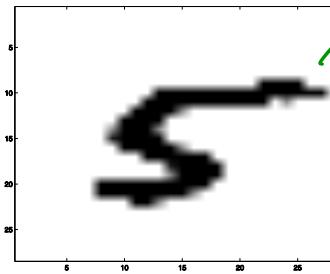
Here  $x_i$  is the  $i$ th coordinate of  $x$ .

# Nearest neighbor

Image to label



Nearest neighbor



Overall:

**error rate = 6%**  
(on test set)

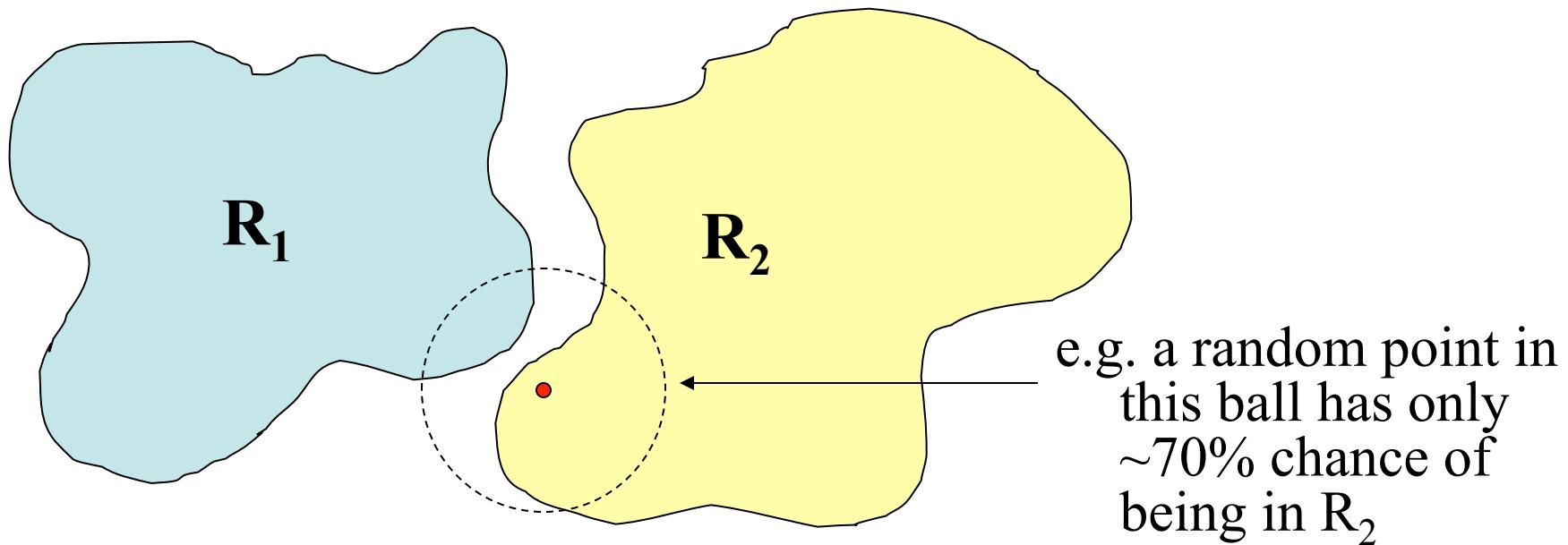
Question: what is the error rate for random guessing?

# What does it get wrong?

Who knows... but here's a hypothesis:

Each digit corresponds to some connected region of  $\mathbb{R}^{784}$ .

Some of the regions come close to each other; problems occur at these boundaries.



# Examples of errors

Query					
NN					

# “Boost” the probability of success

Analogy: suppose a (biased) coin has

$$\Pr(\text{heads}) = 0.70$$

Flip it 11 times and return the majority vote:

$$\Pr(\text{heads}) = 0.92$$

Therefore: to classify  $x$ , find its  $k$  nearest neighbors (in the training set) and return their majority vote

[Large deviation theory: the foundation of machine learning...]

# Time complexity

To classify a new (unlabeled) image:

find its nearest neighbor in the database (training set) and return that label

Given  $n$  training points in  $\mathbb{R}^d$ , what is the time complexity to label a new image,  $q$ ?

For each image,  $x$ , in the database:

$$= O(n)$$

Compute the Euclidean distance between  $x$  and  $q$ .

Return the  $x_{\min}$  that minimizes this distance.

$$d(q, x) = \sqrt{\sum_{i=1}^d (q_i - x_i)^2} = O(d)$$

$$= O(nd)$$

$$d(q, x) = \sqrt{\sum_{i=1}^d (q_i - x_i)^2} = O(d)$$

# (Simple) Nearest neighbor: pros and cons

## Pros

- Simple
- Flexible
- Excellent performance on a wide range of tasks

## Cons

- Algorithmic:
  - Time consuming: with  $n$  training points in  $\mathbb{R}^d$ , time to label a new point is  $O(nd)$ .
- Statistical:
  - This is just memorization, not learning!
  - No insight is gained about the problem.
  - Would prefer a compact classifier.

# Beyond Simple Nearest Neighbor

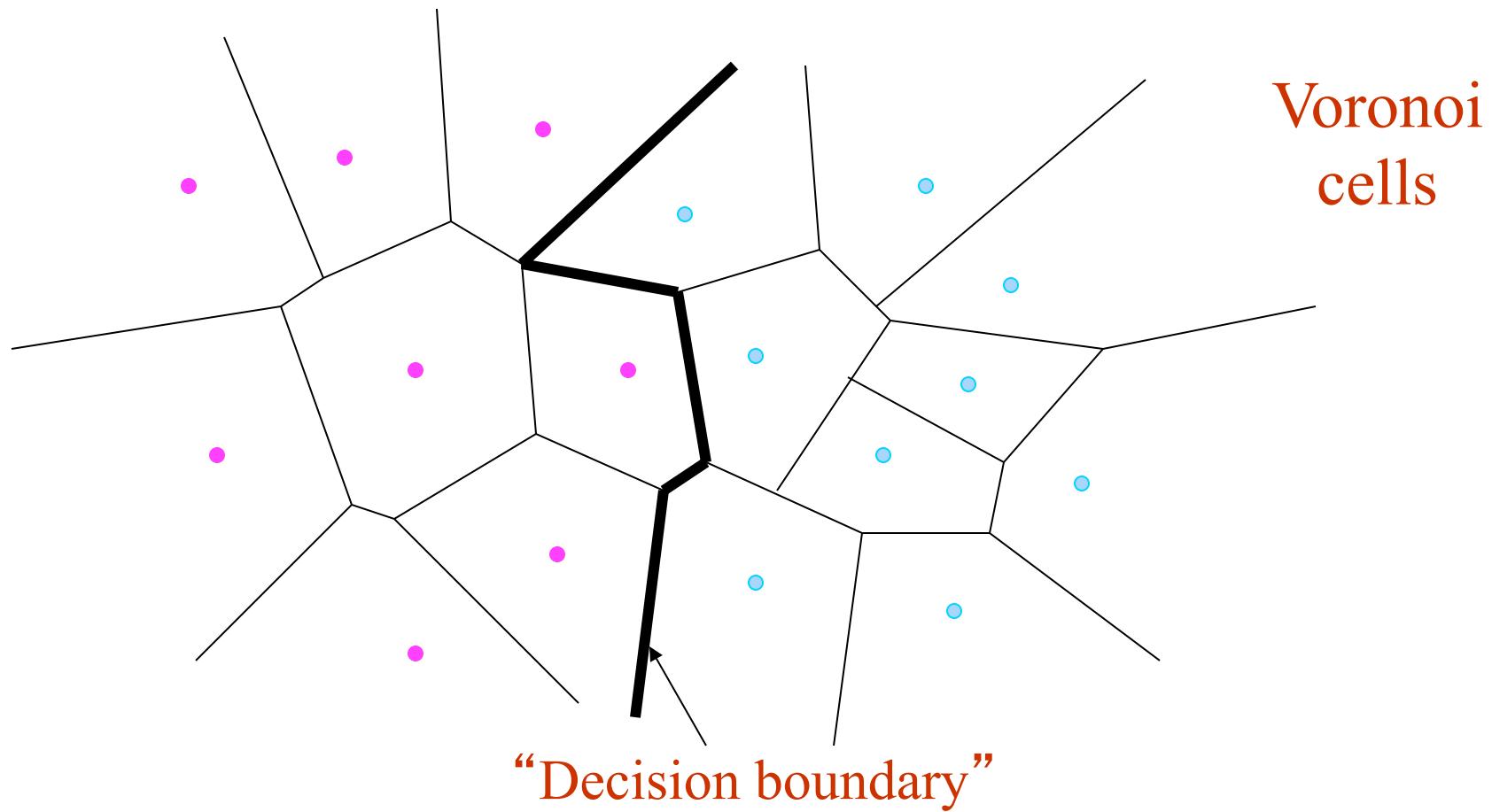
Learning problem:

Design algorithms to perform nearest neighbor classification:

- with **low error**
- and **low memory**: do not store all the training examples!
  - we want **low training set**

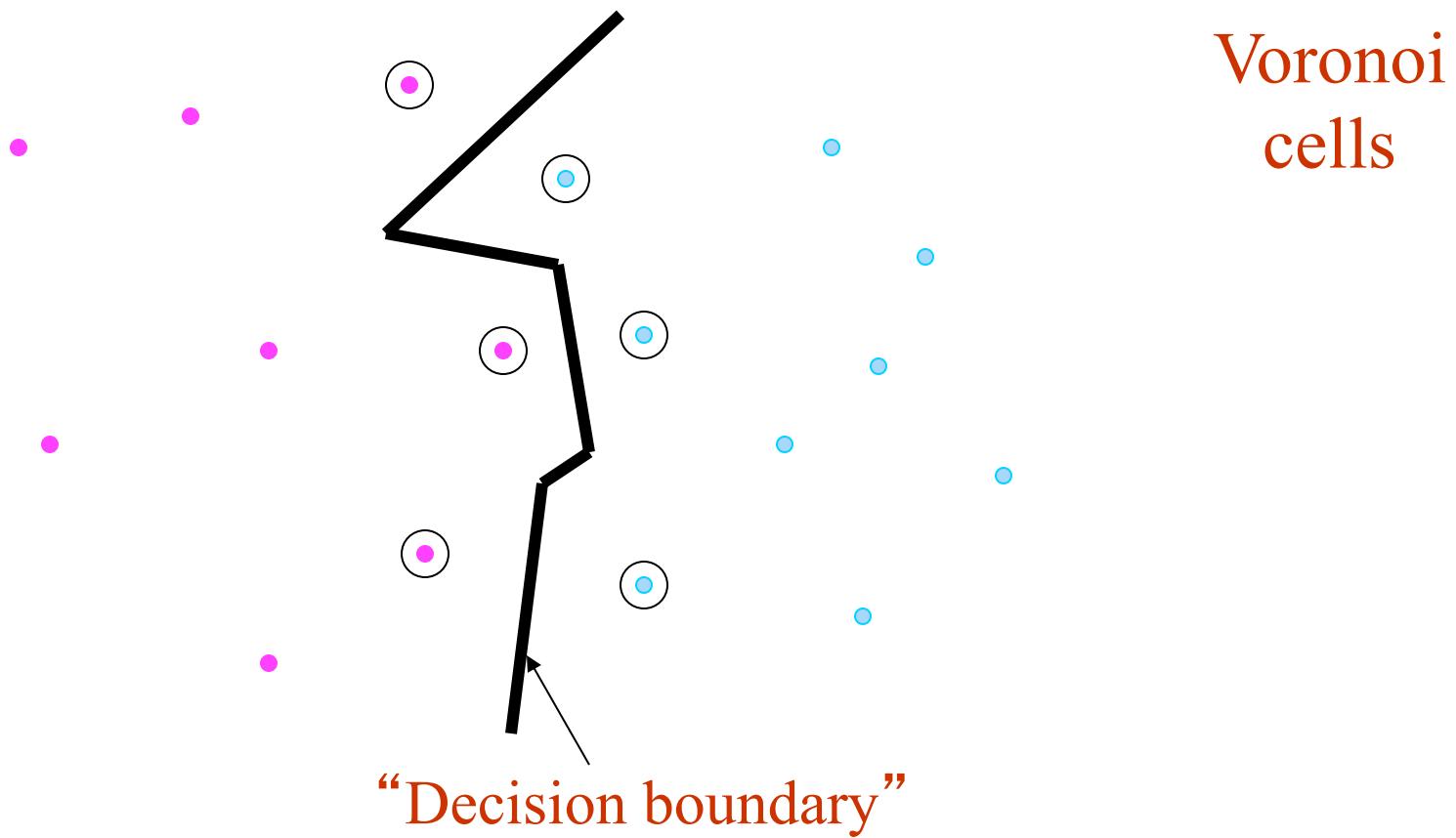
# Prototype selection

A possible fix: instead of the entire training set, just keep a “representative sample”



# Prototype selection

A possible fix: instead of the entire training set, just keep a “representative sample”



# How to pick prototypes?

Idea 1: sample uniformly at random from training data

# prototypes	% error rate
20000	6.0
10000	7.4
5000	8.5
2000	10.7
1000	14.3
500	17.8
250	20.8
100	32.3
50	43.0

it's nice algorithm  
to pick random  
Label dataset.

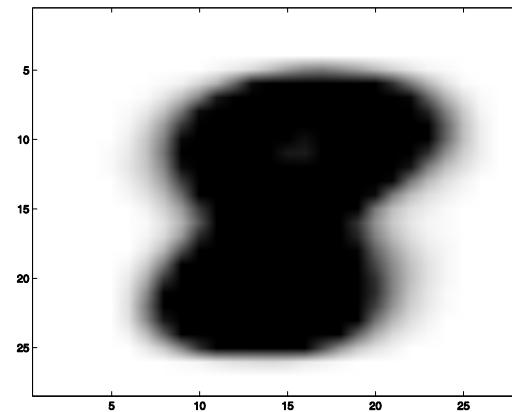
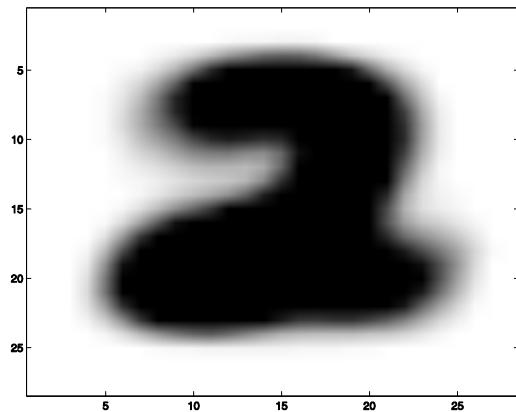
→ wif over 100  
prototypes, you ge f  
32.3 error rate

# How to pick prototypes?

NOTE: They do not have to be actual data points!

Idea 2: one prototype per class: mean of training points  
*has to digit*

Examples:



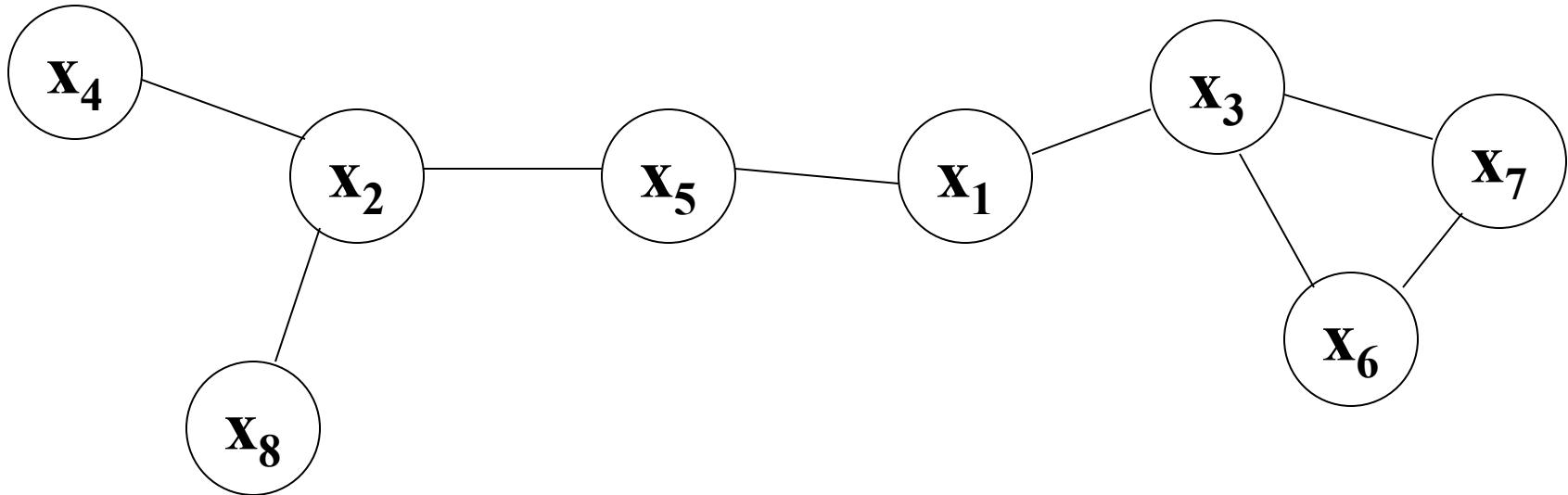
Error = 23%. Number of prototypes =

# Picking prototypes

Other methods? e.g. Gain geometric insight from the neighborhood graph

Node = training point

Edges = connect close neighbors (small distance)



This is an area of machine learning research.

# Generalization

Suppose a set of prototypes performs well on the training data.

==> Will it do well on future test data?

What kinds of classifiers generalize well?

less classifiers generalize well

Answer: the simplest (most compact) classifiers generalize the best. [cf. Occam's razor]

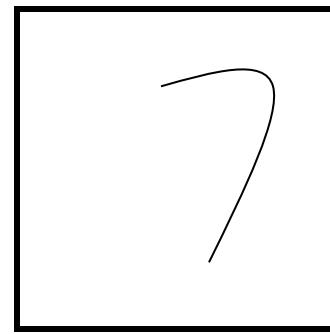
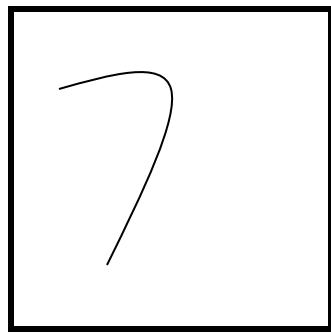
Measure of simplicity:

“Vapnik-Chervonenkis dimension” (later in course)

# Postscript: representation I

The data lie in  $\mathbb{R}^{784}$ , but our particular choice of Euclidean distance was pretty arbitrary.

Also suboptimal, eg.:



... are pretty far apart! (when computing Euclidean distance between vectors of pixels... Can you see why?)

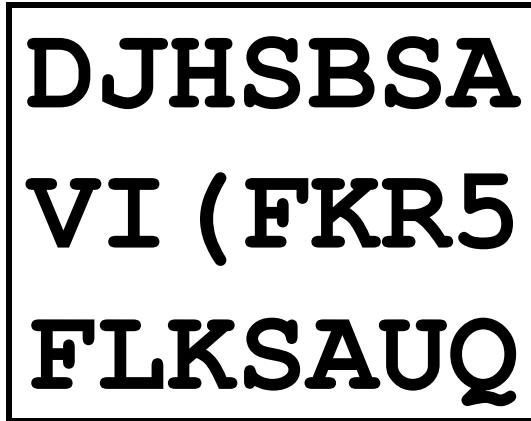
Is there a better distance measure for this application?

# Postscript: representation II

Our 784 features: *pixels, are not very relevant to this task.*

Often data looks like this (*figuratively speaking*):

digit 5:



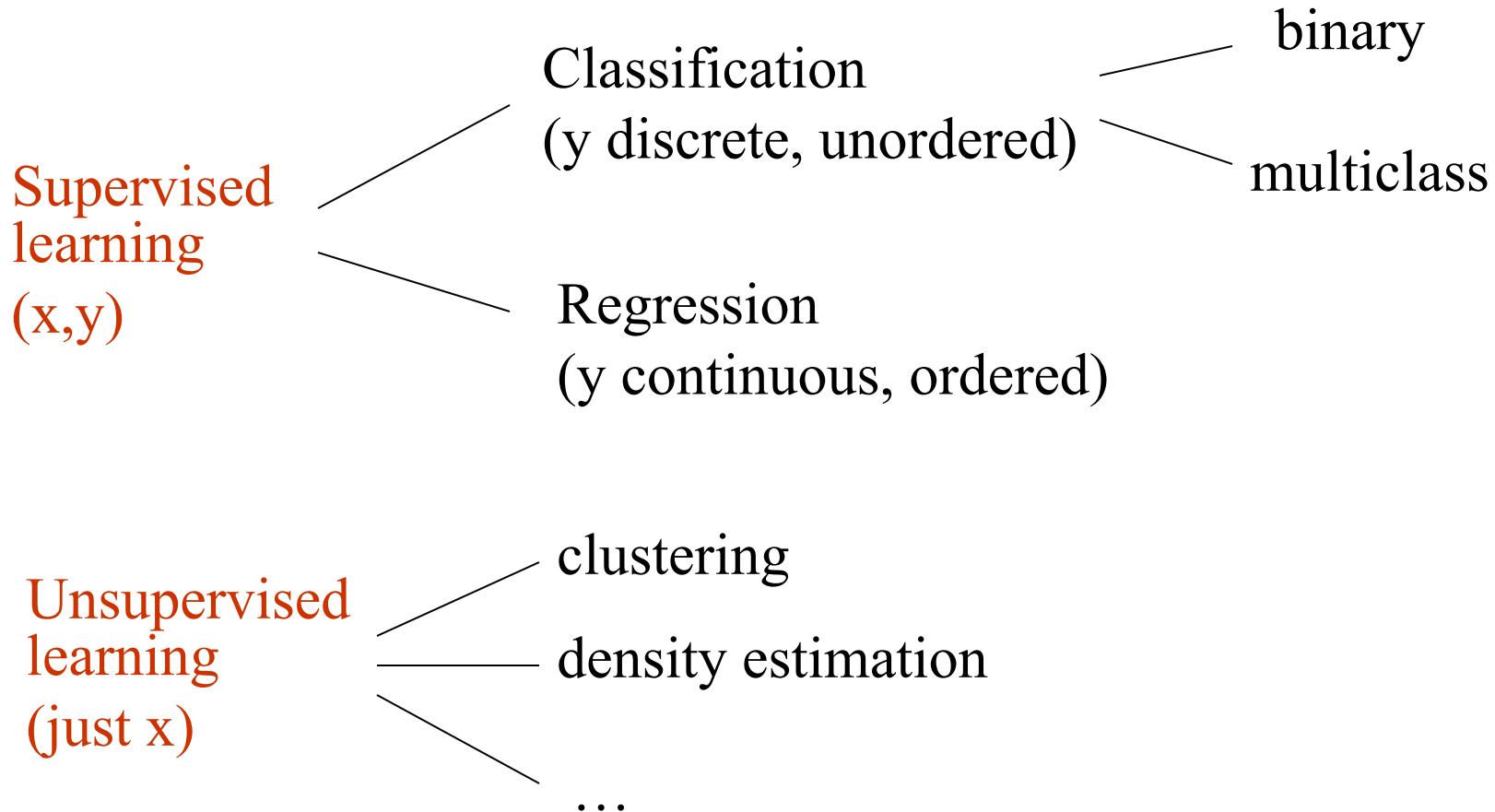
DJHSBSA  
VI (FKR5  
FLKSAUQ

Information buried in a sea of irrelevant features.

Danger for nearest neighbor: match based on irrelevant features.

Feature selection is very important.

# Postscript: learning tasks



All these tasks have a nearest-neighbor solution.