

Text Data Mining:

Predictive and Exploratory Analysis of Text

Jaime Arguello
jarguell@email.unc.edu

August 21, 2013

Introductions

- Hello, my name is _____.
- However, I'd rather be called _____. (optional)
- I'm in the _____ program.
- I'm taking this course because I'd like to learn how to _____.



What is Text Data Mining?

- The science and practice of building and evaluating computer programs that automatically detect or discover interesting and useful things in collections of natural language text

Related Fields

- **Machine Learning:** developing computer programs that improve their performance with “experience”
- **Data Mining:** developing methods that discover patterns within large structured datasets
- **Statistics:** developing methods for the interpretation of data in reaching conclusions with a certain degree of confidence

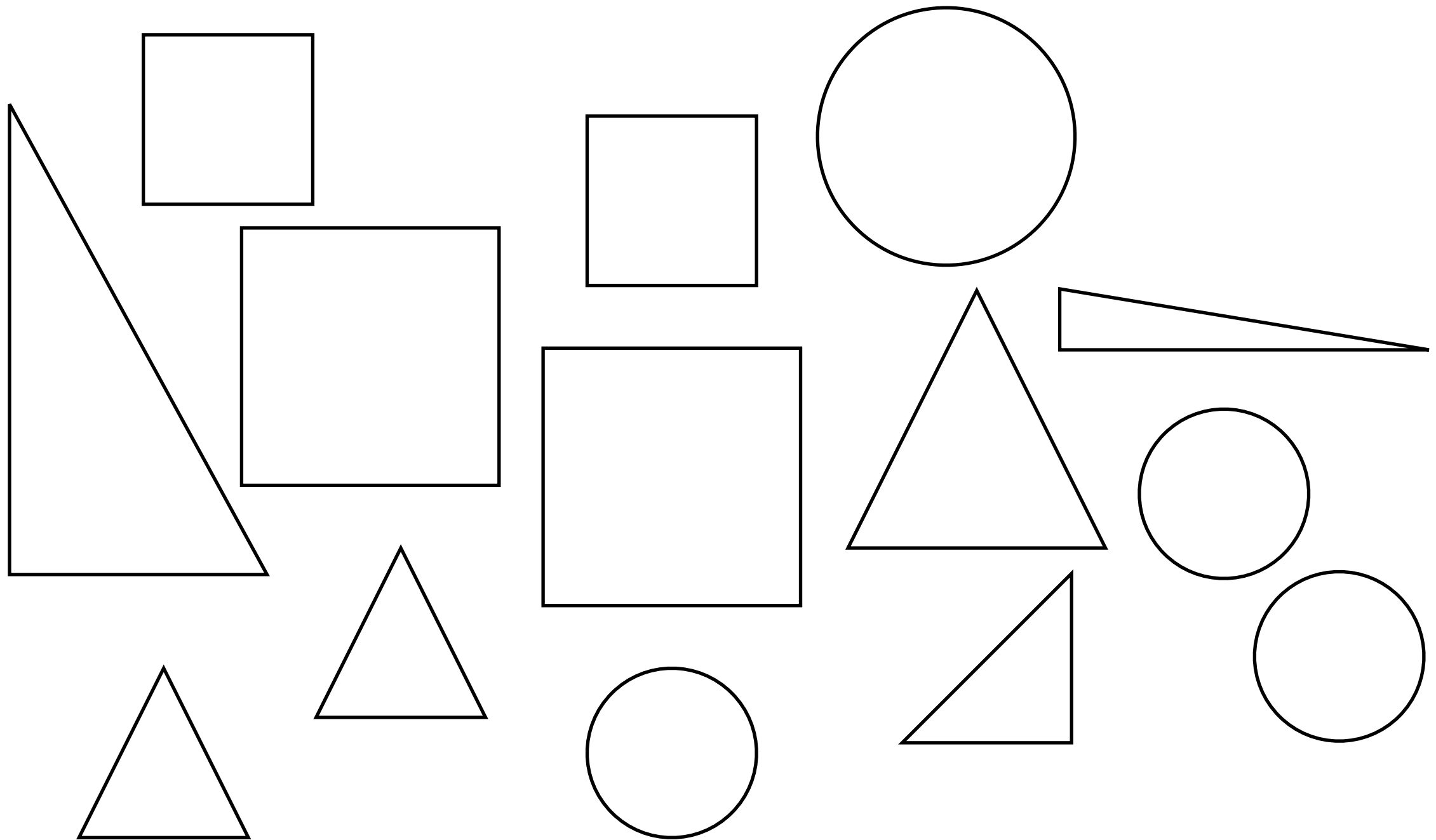
Text Data Mining in this Course

- Predictive Analysis of Text
 - ▶ developing computer programs that automatically recognize or detect a particular concept within a span of text
- Exploratory Analysis of Text:
 - ▶ developing computer programs that automatically discover interesting and useful patterns or trends in text collections

Predictive Analysis: The Big Picture

Predictive Analysis

example: recognizing triangles



Predictive Analysis

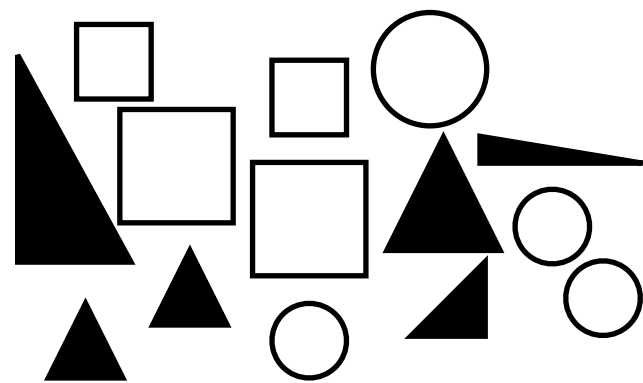
example: recognizing triangles

- We could imagine writing a “triangle detector” by hand:
 - ▶ if shape has three sides, then shape = triangle.
 - ▶ otherwise, shape = other
- Alternatively, we could use supervised machine learning!

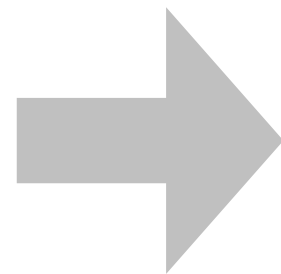
Predictive Analysis

example: recognizing triangles

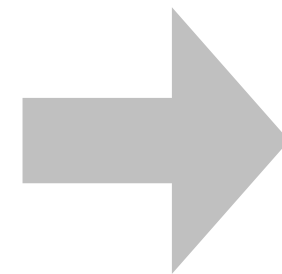
training



labeled examples

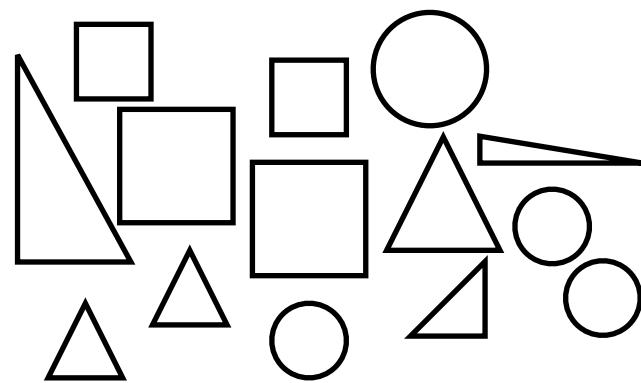


machine
learning
algorithm

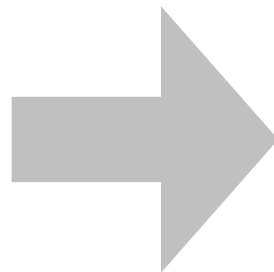


model

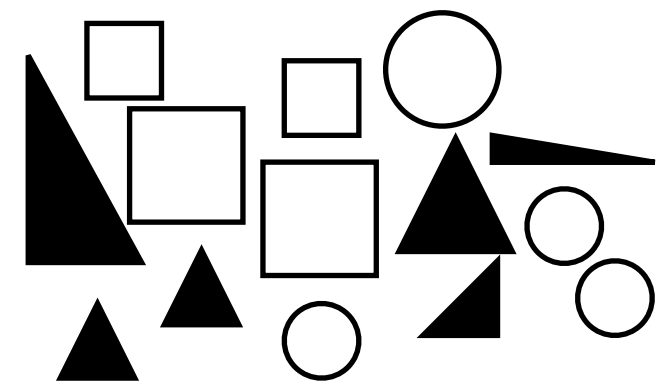
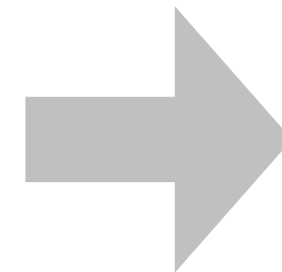
testing



new, unlabeled
examples



model



predictions

Predictive Analysis

basic ingredients

1. **Training data:** a set of examples of the concept we want to automatically recognize
2. **Representation:** a set of features that we believe are useful in recognizing the desired concept
3. **Learning algorithm:** a computer program that uses the training data to learn a predictive model of the concept

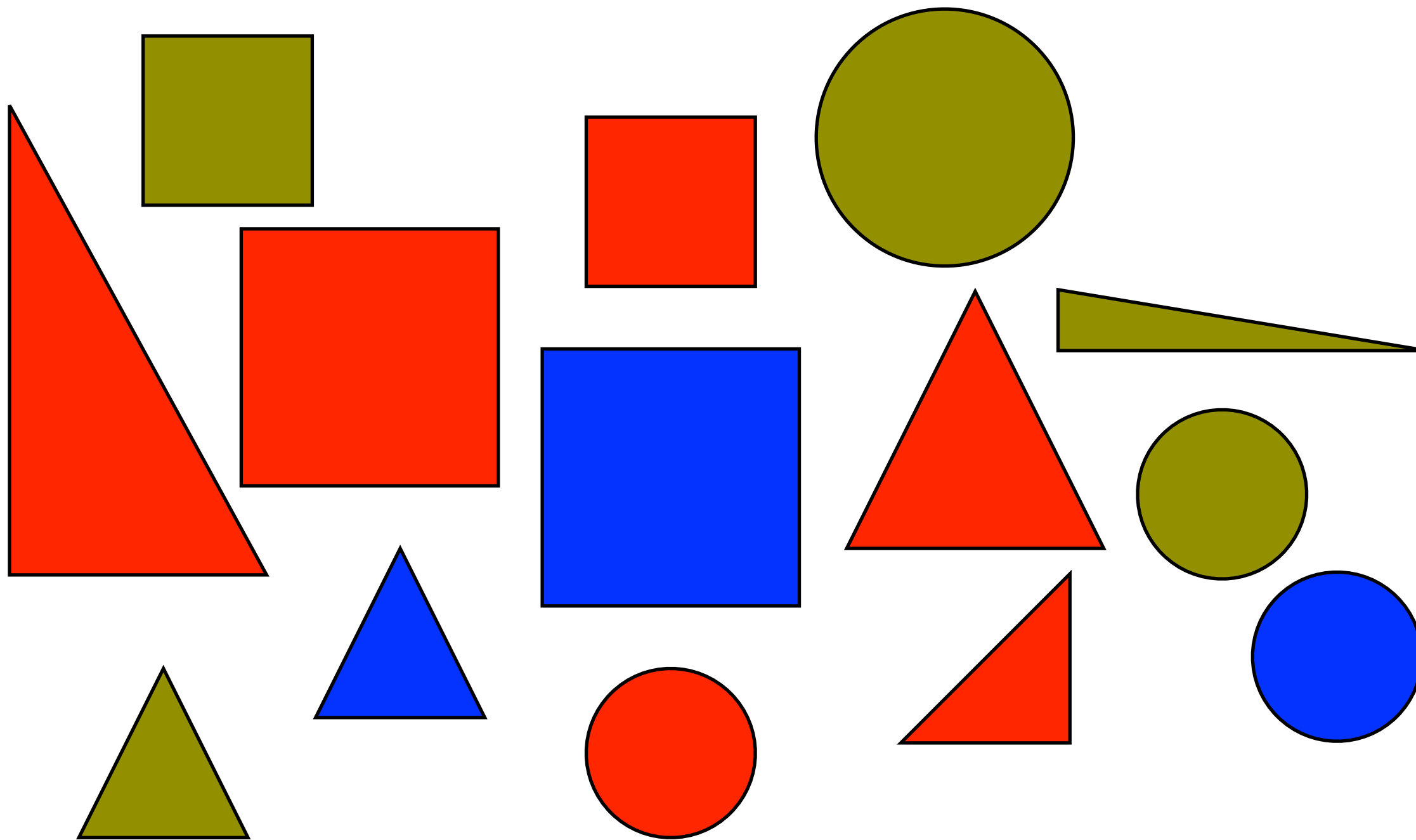
Predictive Analysis

basic ingredients

4. **Model:** a function that describes a predictive relationship between feature values and the presence/absence of the concept
5. **Test data:** a set of previously unseen examples used to estimate the model's effectiveness
6. **Performance metrics:** a set of statistics used measure the predictive effectiveness of the model

Predictive Analysis

representation: features



Predictive Analysis

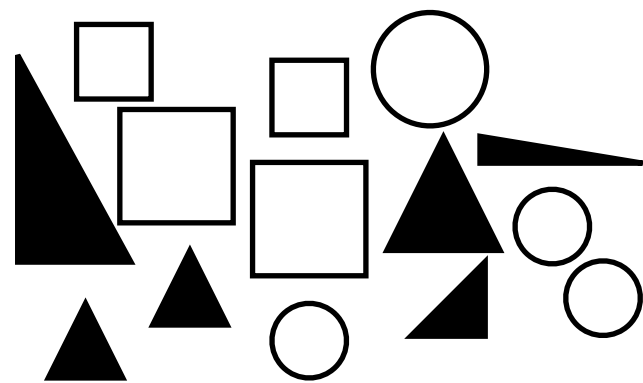
representation: features

color	size	# slides	equal sides	...	label
red	big	3	no	...	yes
green	big	3	yes	...	yes
blue	small	inf	yes	...	no
blue	small	4	yes	...	no
⋮	⋮	⋮	⋮	⋮	⋮
red	big	3	yes	...	yes

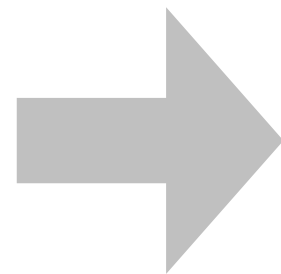
Predictive Analysis

example: recognizing triangles

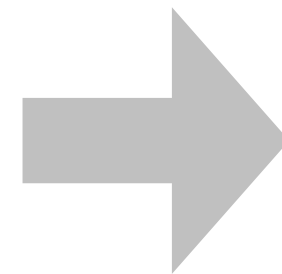
training



labeled examples

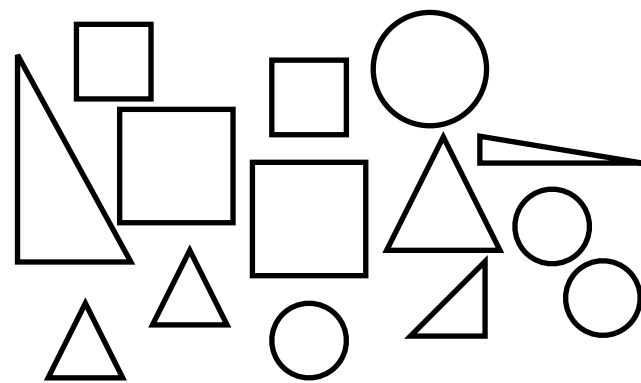


machine
learning
algorithm

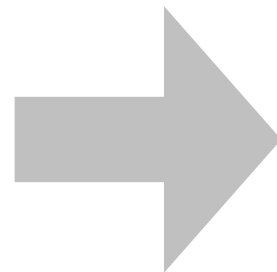


model

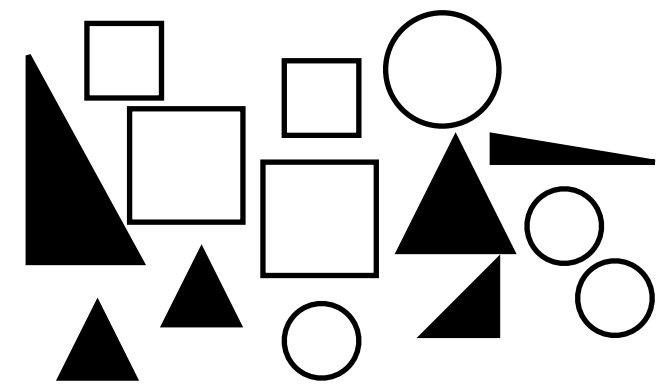
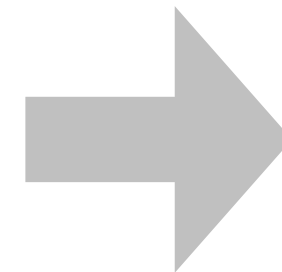
testing



new, unlabeled
examples



model



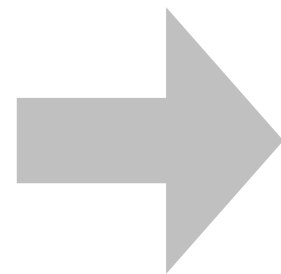
predictions

Predictive Analysis

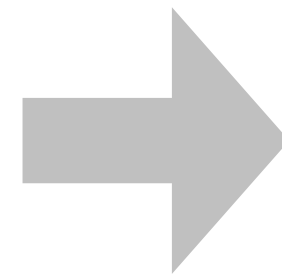
example: recognizing triangles

training

color	size	sides	equal sides	...	label
red	big	3	no	...	yes
green	big	3	yes	...	yes
blue	small	inf	yes	...	no
blue	small	4	yes	...	no
⋮	⋮	⋮	⋮	⋮	⋮
red	big	3	yes	...	yes



machine
learning
algorithm

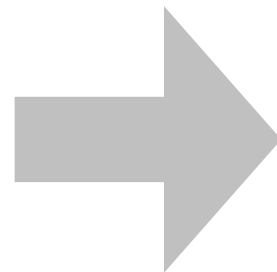


model

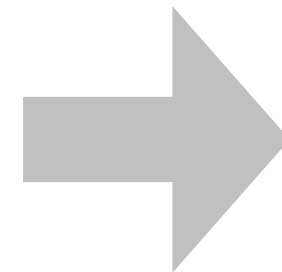
labeled examples

testing

color	size	sides	equal sides	...	label
red	big	3	no	...	???
green	big	3	yes	...	???
blue	small	inf	yes	...	???
blue	small	4	yes	...	???
⋮	⋮	⋮	⋮	⋮	???
red	big	3	yes	...	???



model



color	size	sides	equal sides	...	label
red	big	3	no	...	yes
green	big	3	yes	...	yes
blue	small	inf	yes	...	no
blue	small	4	yes	...	no
⋮	⋮	⋮	⋮	⋮	⋮
red	big	3	yes	...	yes

new, unlabeled
examples

predictions

Predictive Analysis

basic ingredients: the focus in this course

- ▶ 1. **Training data:** a set of examples of the concept we want to automatically recognize
- ▶ 2. **Representation:** a set of features that we believe are useful in recognizing the desired concept
- ▶ 3. **Learning algorithm:** uses the training data to learn a predictive model of the “concept”

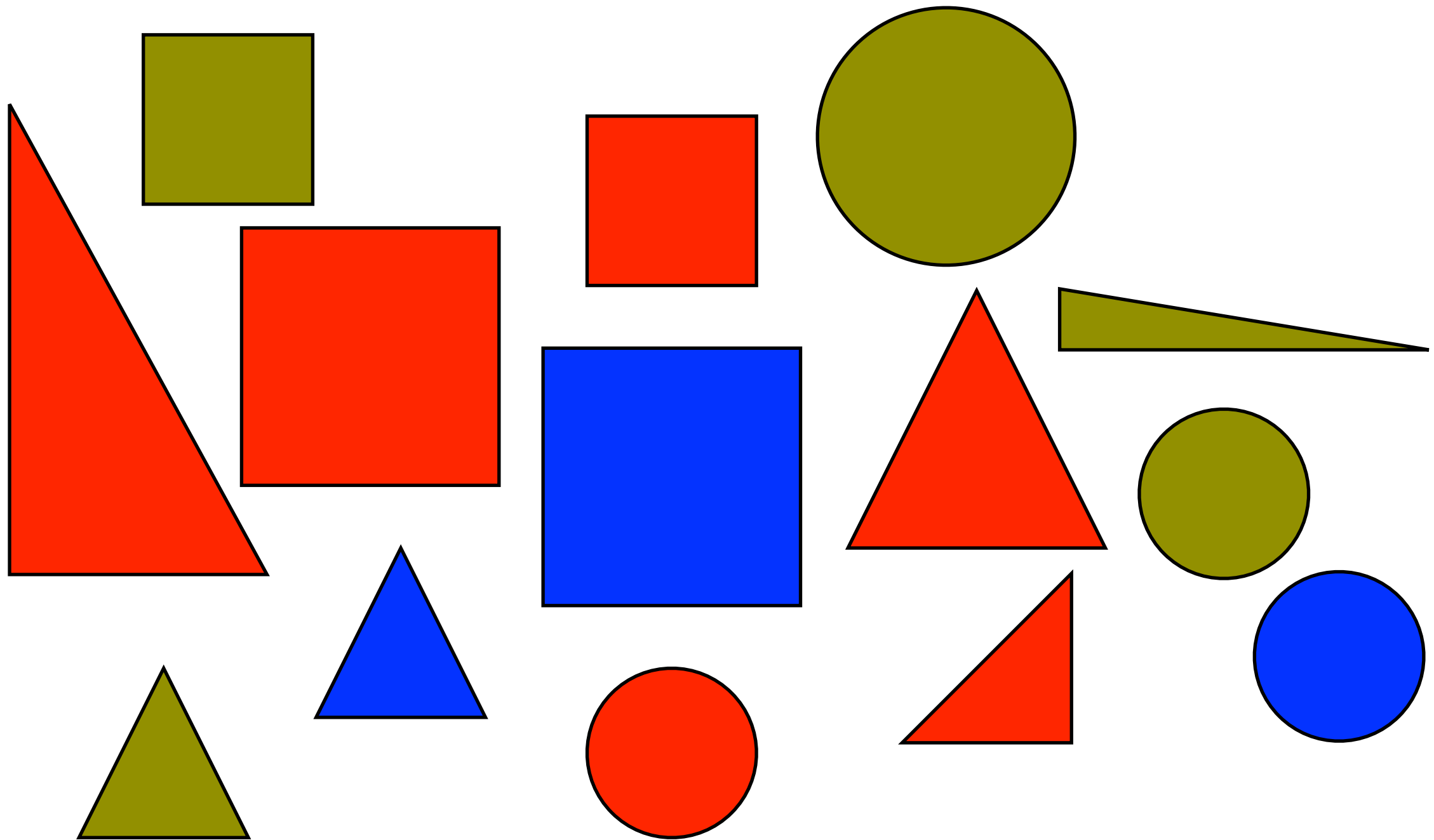
Predictive Analysis

basic ingredients: the focus in this course

- ▶ 3. **Model:** describes a predictive relationship between feature values and the presence/absence of the concept
- ▶ 4. **Test data:** a set of previously unseen examples used to estimate the model's effectiveness
- ▶ 5. **Performance metrics:** a set of statistics used measure the predictive effectiveness of the model

Training data + Representation

what could possibly go wrong?



Training data + Representation

what could possibly go wrong?

color	size	90 deg. angle	equal sides	...	label
red	big	yes	no	...	yes
green	big	no	yes	...	yes
blue	small	no	yes	...	no
blue	small	yes	yes	...	no
⋮	⋮	⋮	⋮	⋮	⋮
red	big	no	yes	...	yes

Training data + Representation

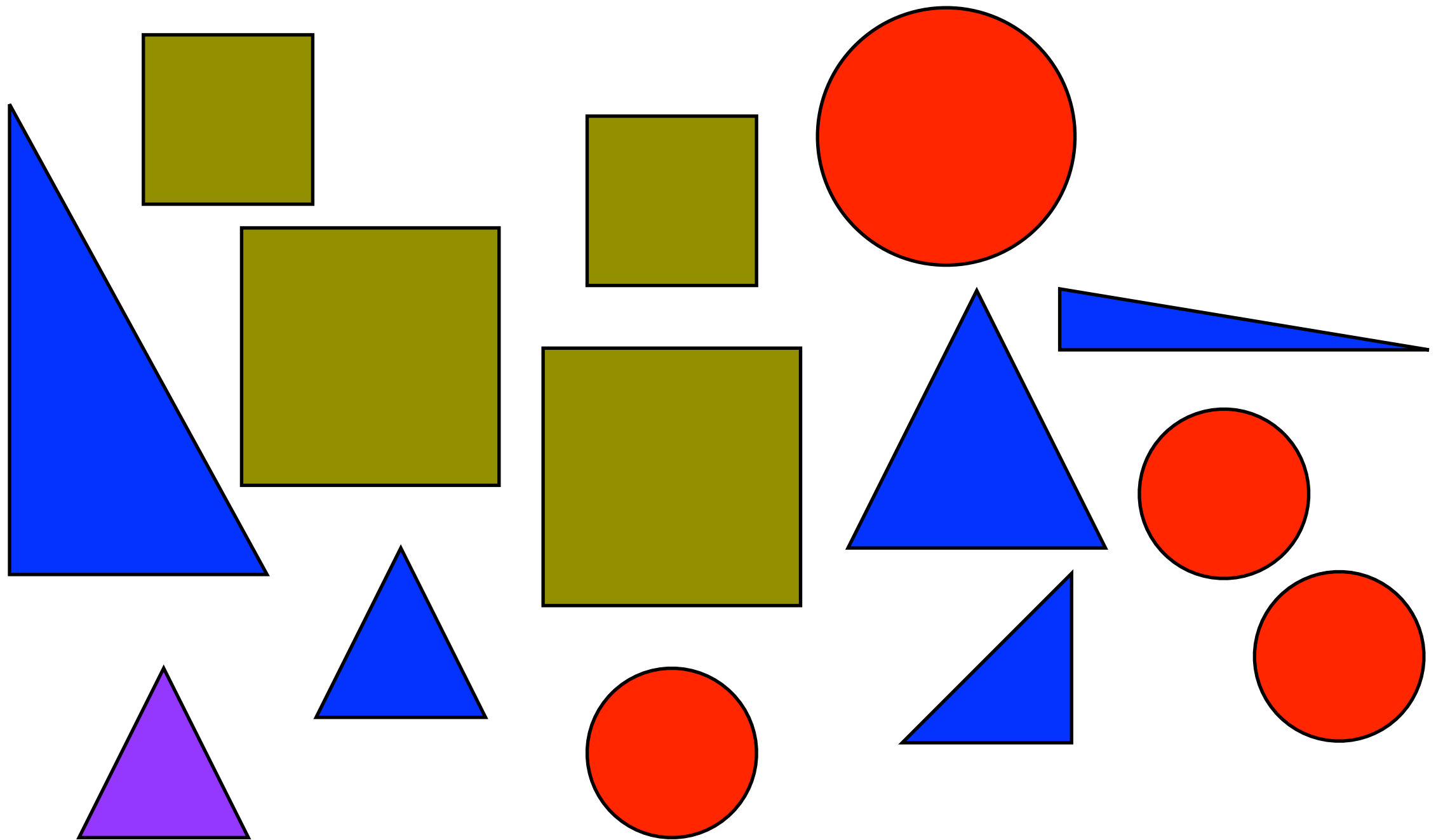
what could possibly go wrong?

color	size	90 deg. angle	equal sides	...	label
red	big	yes	no	...	yes
green	big	no	yes	...	yes
blue	small	no	yes	...	no
blue	small	yes	yes	...	no
⋮	⋮	⋮	⋮	⋮	⋮
red	big	no	yes	...	yes

1. bad feature representation!

Training data + Representation

what could possibly go wrong?



Training data + Representation

what could possibly go wrong?

color	size	# slides	equal sides	...	label
blue	big	3	no	...	yes
blue	big	3	yes	...	yes
red	small	inf	yes	...	no
green	small	4	yes	...	no
⋮	⋮	⋮	⋮	⋮	⋮
blue	big	3	yes	...	yes

Training data + Representation

what could possibly go wrong?

color	size	# slides	equal sides	...	label
blue	big	3	no	...	yes
blue	big	3	yes	...	yes
red	small	inf	yes	...	no
green	small	4	yes	...	no
⋮	⋮	⋮	⋮	⋮	⋮
blue	big	3	yes	...	yes

2. bad data + misleading features

Training data + Representation

what could possibly go wrong?

color	size	# slides	equal sides	...	label
white	big	3	no	...	yes
white	big	3	no	...	no
white	small	inf	yes	...	yes
white	small	4	yes	...	no
⋮	⋮	⋮	⋮	⋮	⋮
white	big	3	yes	...	yes

3. bad training data!

Evaluating a Model

example: recognizing triangles

1. Make predictions on data that is “unlabeled”
 - ▶ we know what the true labels are, but the model doesn’t “see” or use these labels
2. Evaluate based on some metric that compares the model’s predicted labels with the known true labels

color	size	sides	equal sides	...	label
red	big	3	no	...	???
green	big	3	yes	...	???
blue	small	inf	yes	...	???
blue	small	4	yes	...	???
⋮	⋮	⋮	⋮	⋮	???
red	big	3	yes	...	???

new, unlabeled
examples

testing



color	size	sides	equal sides	...	label
red	big	3	no	...	yes
green	big	3	yes	...	yes
blue	small	inf	yes	...	no
blue	small	4	yes	...	no
⋮	⋮	⋮	⋮	⋮	⋮
red	big	3	yes	...	yes

predictions

Evaluating a Model

example: recognizing triangles

- Many evaluation metrics can be generated from a contingency table

		true	
		triangle	other
predicted	triangle	A	B
	other	C	D

- What number(s) do we want to maximize?
- What number(s) do we want to minimize?

Evaluating a Model

example: recognizing triangles

- **True positives (A):** number of triangles correctly predicted as triangles
- **False positives (B):** number of “other” incorrectly predicted as triangles
- **False negatives (C):** number of triangles incorrectly predicted as “other”
- **True negatives (D):** number of “other” correctly predicted as “other”

		true	
		triangle	other
predicted	triangle	A	B
	other	C	D

Evaluating a Model

your first metric: accuracy

- **Accuracy:** percentage of predictions that are correct (i.e., true positives and true negatives)

$$(\text{?} + \text{?})$$

$$(\text{?} + \text{?} + \text{?} + \text{?})$$

true

		true	
		triangle	other
predicted	triangle	A	B
	other	C	D

Evaluating a Model

your first metric: accuracy

- **Accuracy:** percentage of predictions that are correct (i.e., true positives and true negatives)

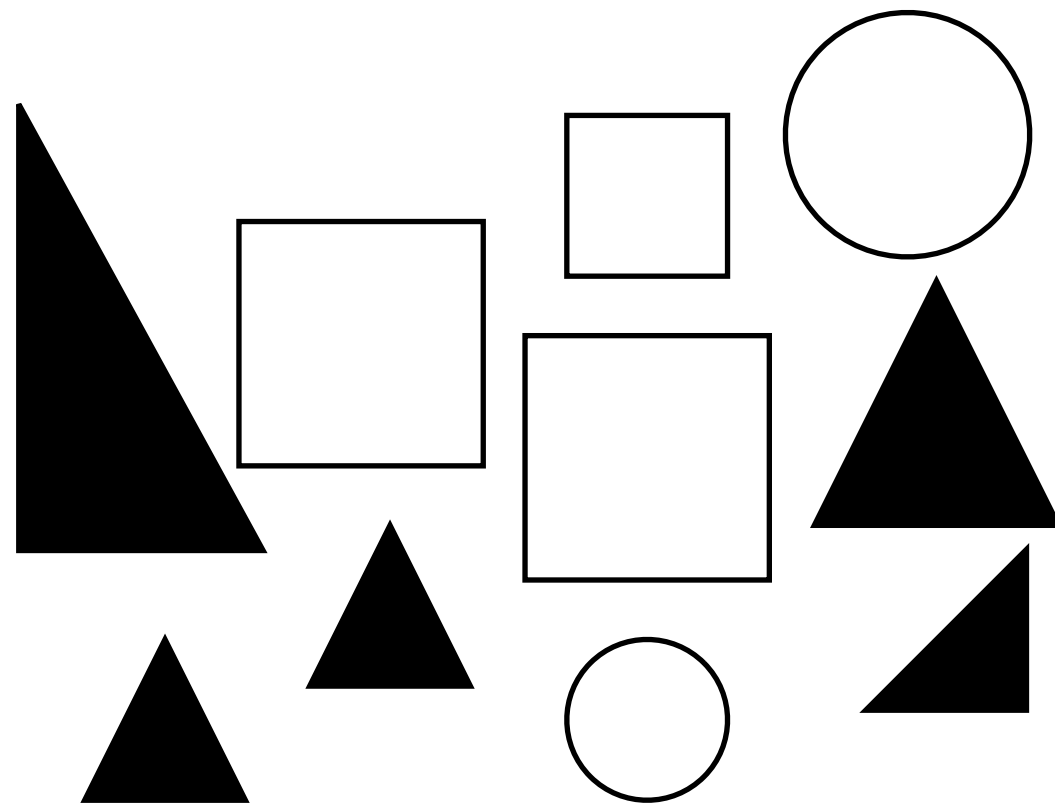
$$\frac{(A + D)}{(A + B + C + D)}$$

		true	
		triangle	other
predicted	triangle	A	B
	other	C	D

Evaluating a Model

your first metric: accuracy

- **Accuracy:** percentage of predictions that are correct (i.e., true positives and true negatives)

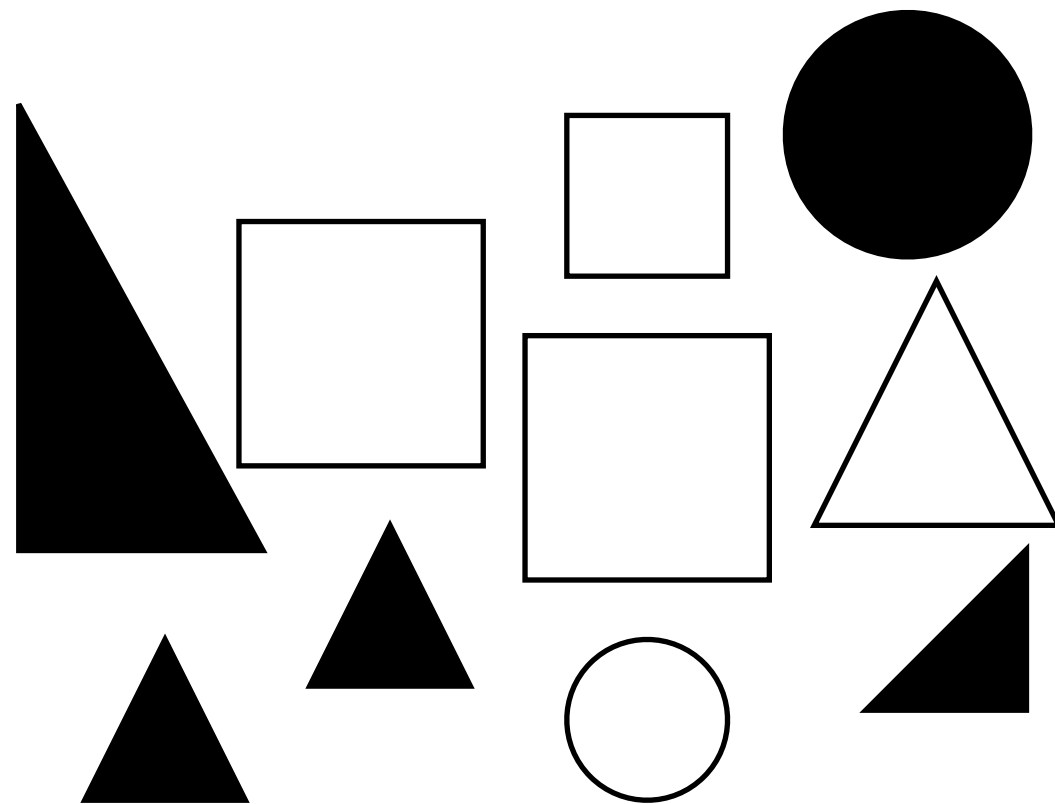


- What is the accuracy of this model?

Evaluating a Model

your first metric: accuracy

- **Accuracy:** percentage of predictions that are correct (i.e., true positives and true negatives)

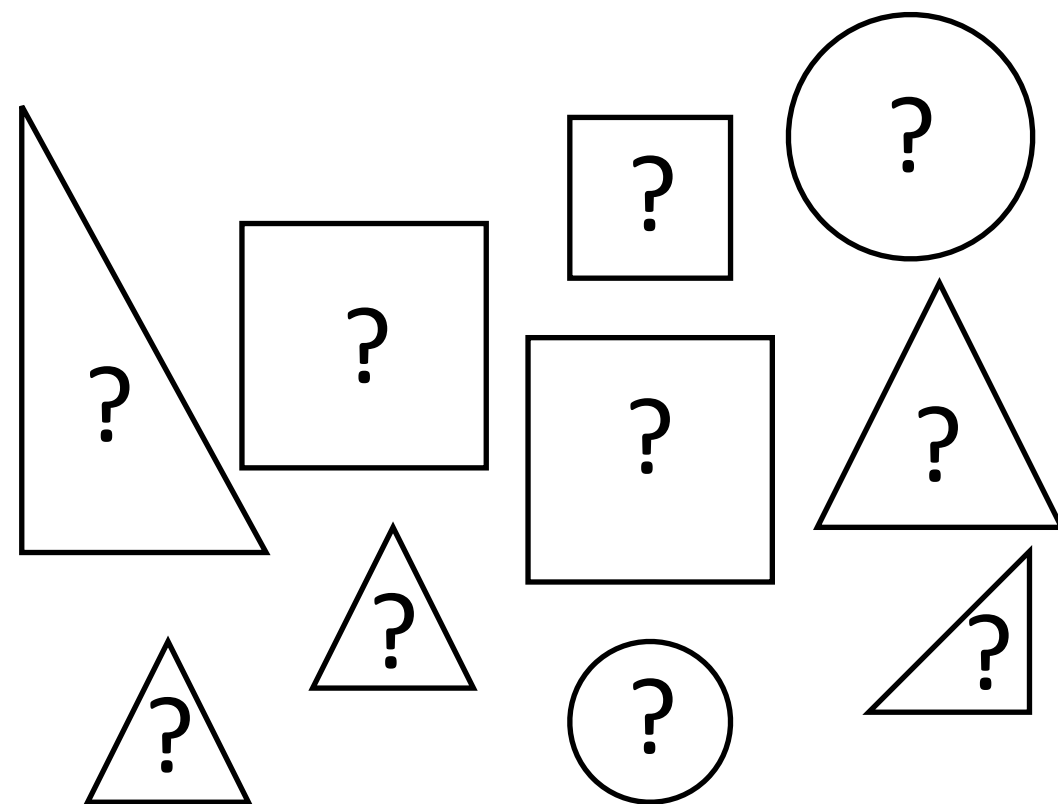


- What is the accuracy of this model?

Evaluating a Model

what could possibly go wrong?

- Interpreting the value of a metric on a particular data set requires some thinking ...

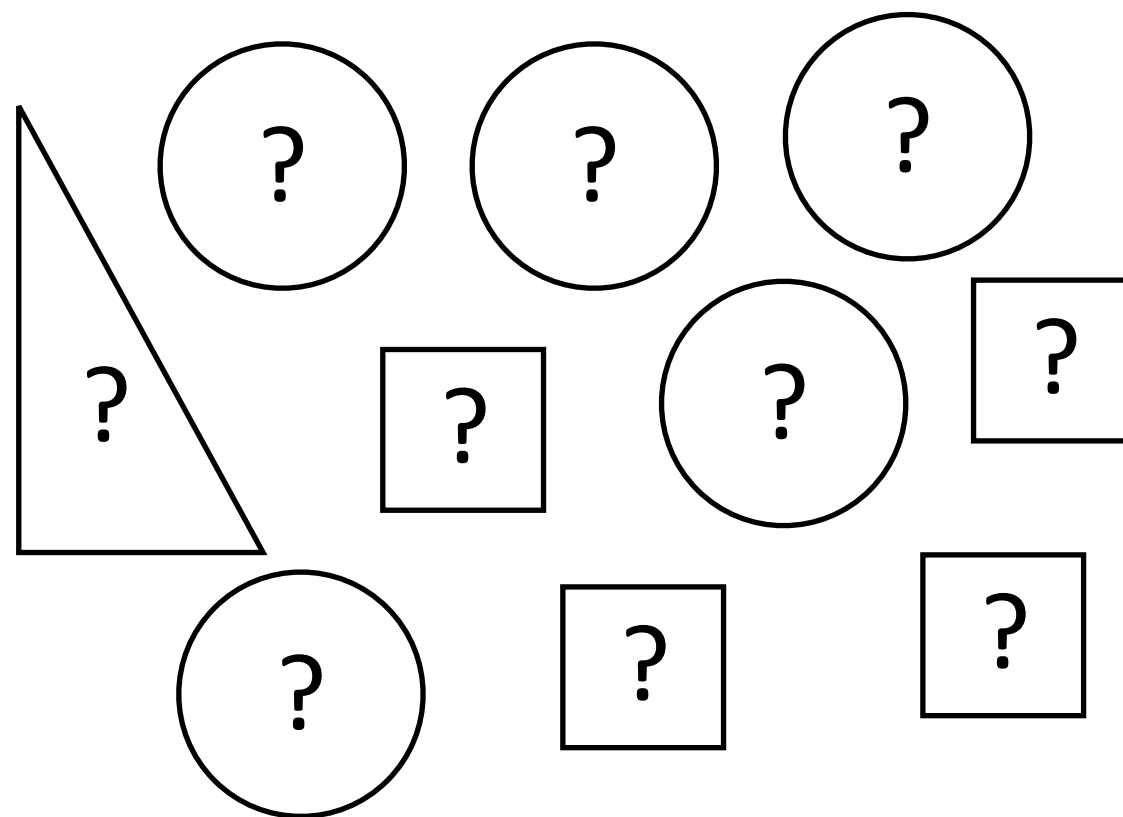


- On this dataset, what would be the expected accuracy of a model that does NO learning (degenerate baseline)?

Evaluating a Model

what could possibly go wrong?

- Interpreting the value of a metric on a particular data set requires some thinking ...



- On this dataset, what would be the expected accuracy of a model that does NO learning (degenerate baseline)?

Text Data Mining:

Predictive and Exploratory Analysis of Text

Jaime Arguello
jarguell@email.unc.edu

August 26, 2013

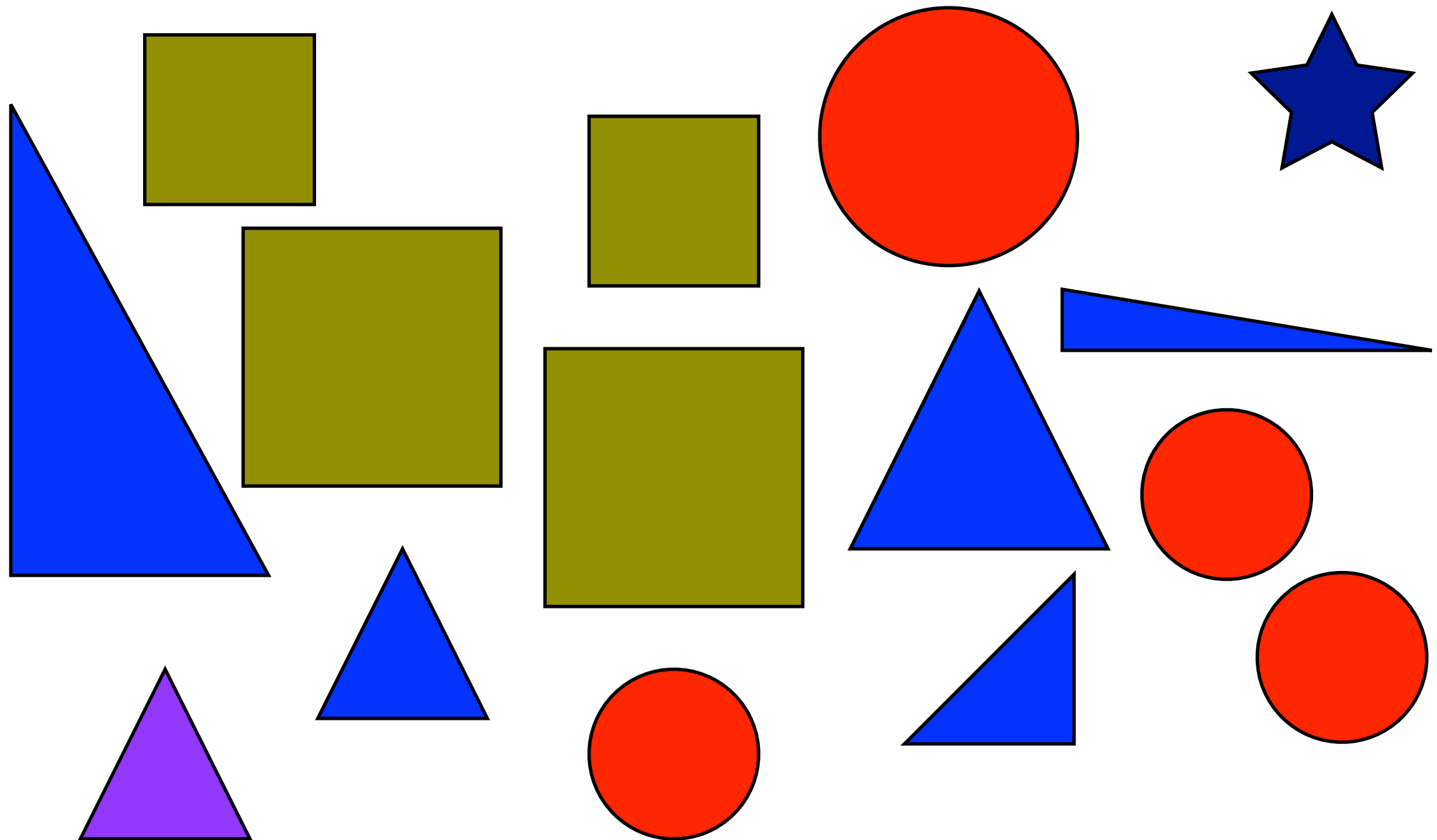
Text Data Mining in this Course

- Predictive Analysis of Text
 - ▶ developing computer programs that automatically recognize a particular concept within a span of text
- Exploratory Analysis of Text:
 - ▶ developing computer programs that automatically discover useful patterns or trends in text collections

Exploratory Analysis: The Big Picture

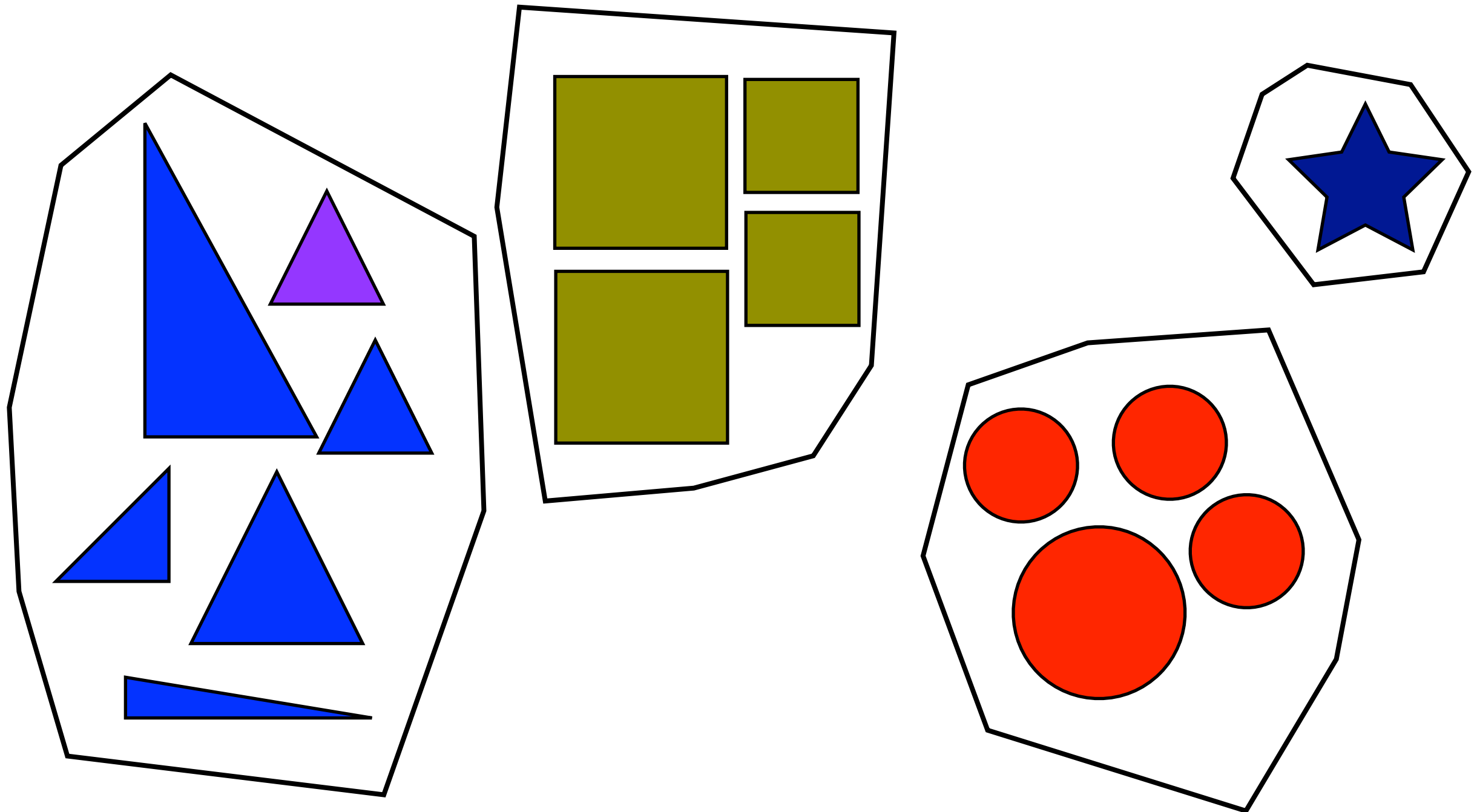
Exploratory Analysis

example: clustering shapes



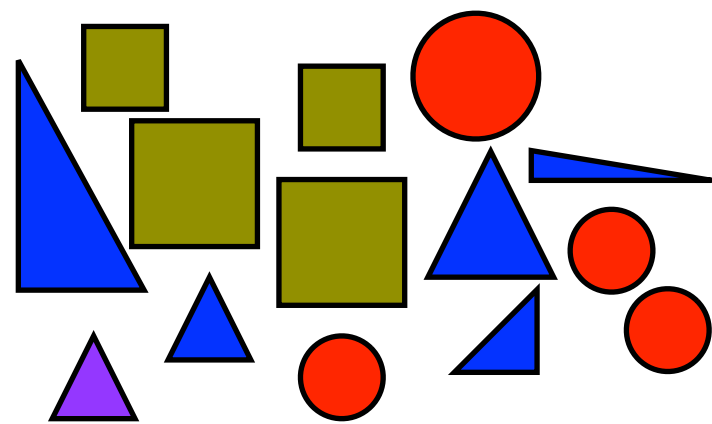
Exploratory Analysis

example: clustering shapes

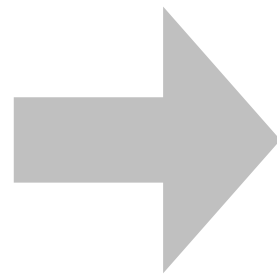


Exploratory Analysis

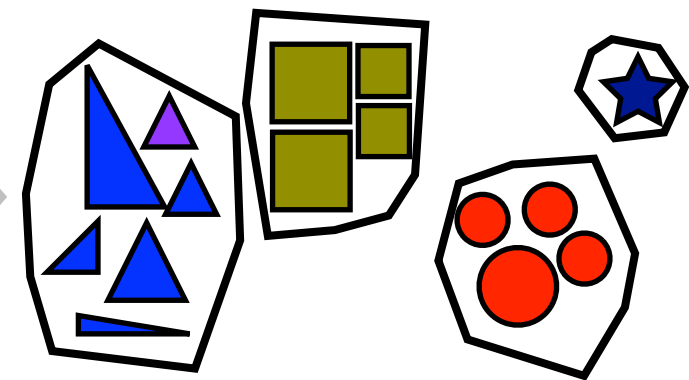
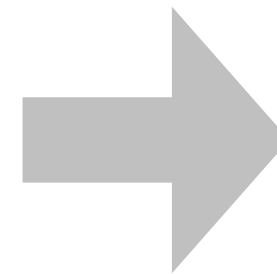
example: clustering shapes



unlabeled
examples



clustering
algorithm



Exploratory Analysis

basic ingredients

1. **Data:** a set of examples that we want to automatically analyze in order to discover interesting trends
2. **Representation:** a set of features that we believe are useful in describing the data (i.e., its main attributes)
3. **Similarity Metric:** a measure of similarity between two examples that is based on their feature values
4. **Clustering algorithm:** an algorithm that assigns items with similar feature values to the same group

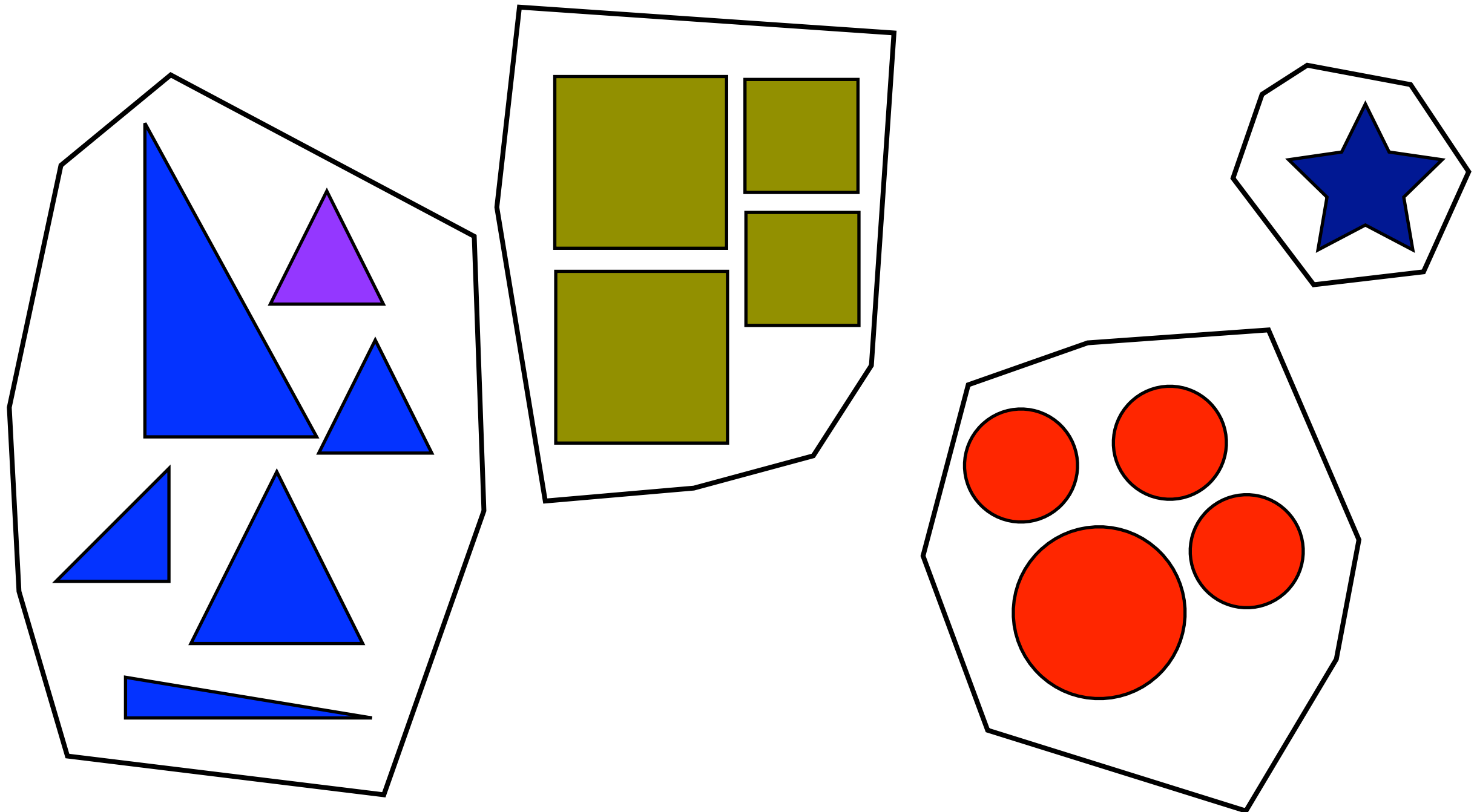
Exploratory Analysis

representation: features

color	size	# slides	equal sides	...	shape
blue	big	3	no	...	triangle
blue	big	3	yes	...	triangle
red	small	inf	yes	...	circle
green	small	4	yes	...	square
⋮	⋮	⋮	⋮	⋮	⋮
blue	big	3	yes	...	triangle

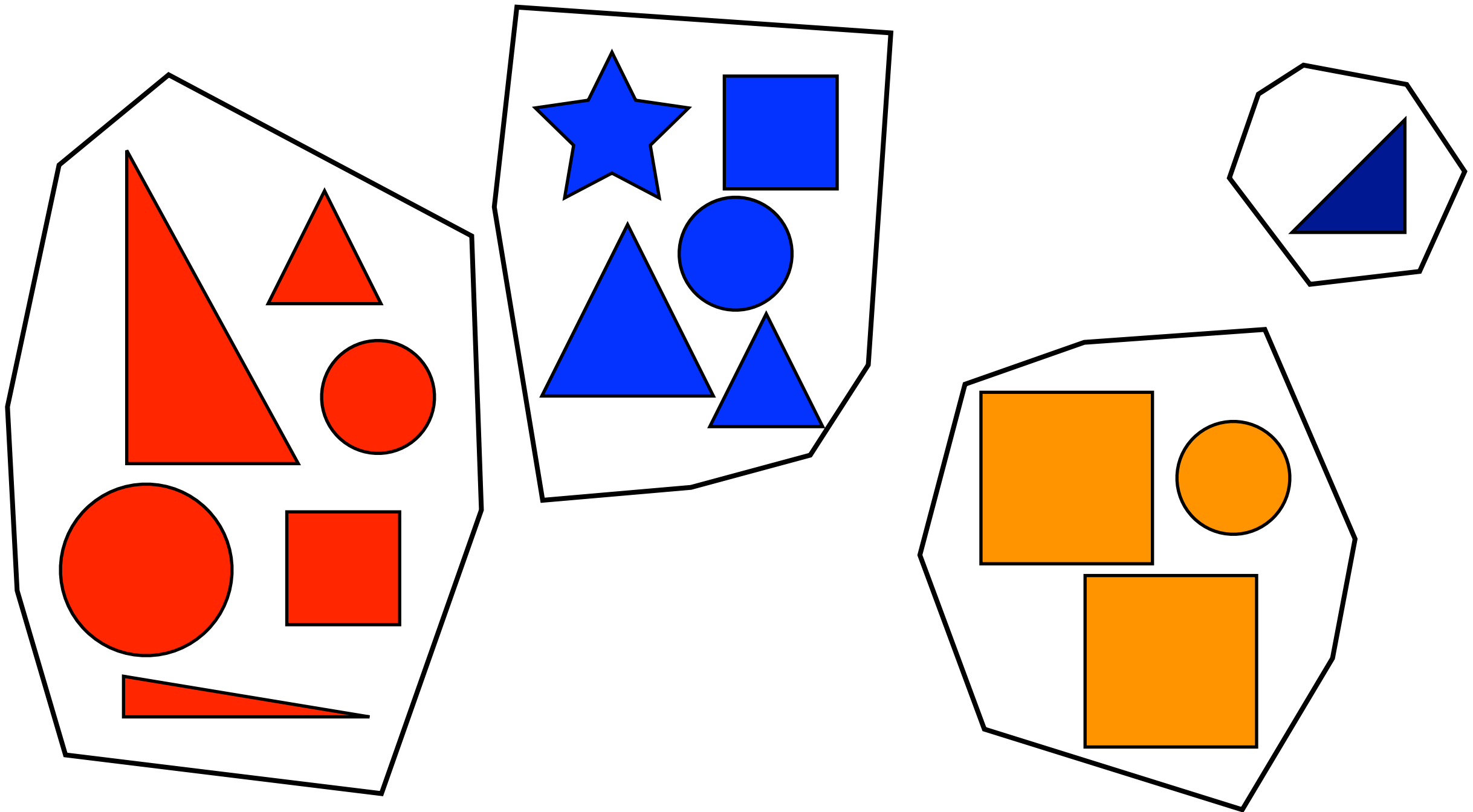
Representation

what could possibly go wrong?






Representation

what could possibly go wrong?



Exploratory Analysis

basic ingredients: the focus in this course

1. **Data:** a set of examples that we want to automatically analyze in order to discover interesting trends
-  2. **Representation:** a set of features that we believe are useful in describing the data
-  3. **Similarity Metric:** a measure of similarity between two examples that is based on their feature values
-  4. **Clustering algorithm:** an algorithm that assigns items with similar feature values to the same group

Text Data Mining in this Course

- Predictive Analysis of Text
 - ▶ developing computer programs that automatically recognize or detect a particular concept within a span of text
- Exploratory Analysis of Text:
 - ▶ developing computer programs that automatically discover interesting and useful patterns or trends in text collections

Predictive Analysis of Text

examples we'll cover in class

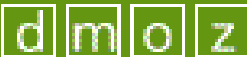
- Topic Categorization
- Opinion Mining
- Sentiment/Affect Analysis
- Bias Detection
- Information Extraction and Relation Learning
- Text-driven Forecasting
- Temporal Summarization

Predictive Analysis of Text

example applications

- **Topic Categorization:** automatically assigning documents to a set of pre-defined topical categories

Topic Categorization

 open directory project

In partnership with
Aol Search.

[about dmoz](#) | [dmoz blog](#) | [suggest URL](#) | [help](#) | [link](#) | [editor login](#)

[advanced](#)

Arts
[Movies](#), [Television](#), [Music](#)...

Games
[Video Games](#), [RPGs](#), [Gambling](#)...

Kids and Teens
[Arts](#), [School Time](#), [Teen Life](#)...

Reference
[Maps](#), [Education](#), [Libraries](#)...

Shopping
[Clothing](#), [Food](#), [Gifts](#)...

World
[Català](#), [Dansk](#), [Deutsch](#), [Español](#), [Français](#), [Italiano](#), [日本語](#), [Nederlands](#), [Polski](#), [Русский](#), [Svenska](#)...

Business
[Jobs](#), [Real Estate](#), [Investing](#)...

Health
[Fitness](#), [Medicine](#), [Alternative](#)...

News
[Media](#), [Newspapers](#), [Weather](#)...

Regional
[US](#), [Canada](#), [UK](#), [Europe](#)...

Society
[People](#), [Religion](#), [Issues](#)...

Computers
[Internet](#), [Software](#), [Hardware](#)...

Home
[Family](#), [Consumers](#), [Cooking](#)...


Recreation
[Travel](#), [Food](#), [Outdoors](#), [Humor](#)...

Science
[Biology](#), [Psychology](#), [Physics](#)...

Sports
[Baseball](#), [Soccer](#), [Basketball](#)...

[Become an Editor](#)

Help build the largest human-edited directory of the web



Copyright © 2012 Netscape

Topic Categorization



In partnership with
Aol Search.

about dmoz	dmoz blog	suggest URL	help	link	editor login
----------------------------	---------------------------	-----------------------------	----------------------	----------------------	------------------------------

--	--

[Search](#) [advanced](#)

Arts

[Movies](#) [Television](#) [Music...](#)

Games

Video Games, RPGs, Gambling...

Kids and Teens

Arts, School Time, Teen Life...

Reference

[Maps, Education, Libraries...](#)

Shopping

[Clothing, Food, Gifts...](#)

World

[Català](#), [Dansk](#), [Deutsch](#), [Español](#), [Français](#), [Italiano](#), [日本語](#), [Nederlands](#), [Polski](#), [Русский](#), [Svenska](#)...

Business

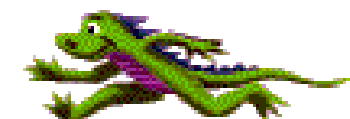
[Jobs, Real Estate, Investing...](#)

Computers

Internet, Software, Hardware...



Become an Editor Help build the largest human-edited directory of the web



Copyright © 2012 Netscape

Predictive Analysis of Text

example applications

- **Opinion Mining:** automatically detecting whether a span of opinionated text expresses a **positive** or **negative** opinion about the item being judged

Opinion Mining

movie reviews

- “Great movie! It kept me on the edge of my seat the whole time. I IMAX-ed it and have no regrets.” positive
- “Waste of time! It sucked!” negative
- “This film should be brilliant. It sounds like a great plot, the actors are first grade, and the supporting cast is good as well, and Stallone is attempting to deliver a good performance. However, it can’t hold up.” negative
- “Trust me, this movie is a masterpiece after you’ve seen it 4+ times.” ???

Predictive Analysis of Text

example applications

- **Sentiment/Effect Analysis:** automatically detecting the emotional state of the author of a span of text (usually from a set of pre-defined emotional states).

Sentiment Analysis

support group posts

- “[I] also found out that the radiologist is doing the biopsy, not a breast surgeon. I am more scared now than when I ...”
fear
- “... My radiologist ‘assured’ me my scan was NOT going to be cancer...she was wrong.”
despair
- “ ... My radiologist did my core biopsy. Not a problem and he did a super job of it.”
hope
- “It's pretty standard for the radiologist to do the biopsy so I wouldn't be concerned on that score.”
hope

Predictive Analysis of Text

example applications

- **Bias detection:** automatically detecting whether the author of a span of text favors a particular viewpoint (usually from a set of pre-defined viewpoints)

Bias Detection

- “Coming [up] next, drug addicted pregnant women no longer have anything to fear from the authorities thanks to the Supreme Court. Both sides on this in a moment.” -- Bill O'Reilly
pro-policy
(vs. anti-policy)
- “Nationalizing businesses, nationalizing banks, is not a solution for the democratic party, it's the objective.” -- Rush Limbaugh
conservative
(vs. liberal)
- “If you're keeping score at home, so far our war in Iraq has created a police state in that country and socialism in Spain. So, no democracies yet, but we're really getting close.” -- Jon Stewart
against war in Iraq
(vs. in favor of)

Predictive Analysis of Text

example applications

- **Information extraction:** automatically detecting that a short sequence of words belongs to (or is an instance of) a particular entity type, for example:
 - ▶ **Person(X)**
 - ▶ **Location(X)**
 - ▶ **TennisPlayer(X)**
 - ▶ ...

Information Extraction

demo

<http://www.boowa.com/>

Predictive Analysis of Text

example applications

- **Relation Learning:** automatically detecting pairs of entities that share a particular relation, for example:
 - ▶ **CEO**(<person>,<company>)
 - ▶ **Capital**(<city>,<country>)
 - ▶ **Mother**(<person>,<person>)
 - ▶ **ConvictedFelon**(<person>,<crime>)
 - ▶ ...

Relation Learning

CEO(<person>, <company>)



[Know Yahoo's Marissa Mayer in 11 facts - CNN.com](#)
www.cnn.com/2012/07/17/...marissa-mayer/index.html



by John D. Sutter - in 846,411 Google+ circles - More by John D. Sutter
Jul 19, 2012 – Here's a quick guide to some of the most interesting and water-cooler-worthy facts about **Marissa Mayer**, who was named CEO of **Yahoo** on
...

<person>, who was named CEO of <company>

Relation Learning

CEO(<person>,<company>)

",who was named CEO of"



[DailyTech - Fisker Appoints New CEO, Eliminates Battery/Engine ...](#)

www.dailytech.com/article.aspx?newsid=25412

4 days ago – Tom LaSorda, **who was named CEO of** Fisker back in February 2012 when founder Henrik Fisker stepped down, is leaving the company, but ...

CEO(Tom LaSorda, Fisker)

[who was named CEO of Yahoo on Monday. Christian Science Monitor](#)

gtp123.com/.../who-was-named-ceo-of-yahoo-on-monday-christian-...

Jul 17, 2012 – You are browsing the archive for **who was named CEO of** Yahoo on Monday. Christian Science Monitor. Avatar of Garland E. Harris ...

[CEO of renamed Sara Lee meat biz chooses Winnetka - Residential ...](#)

www.chicagorealestatedaily.com › Home › Residential News

Aug 7, 2012 – Sean Connolly, **who was named CEO of** Hillshire Brands Co. in January, declines to comment through a company spokesman. Records show ...

CEO(Sean Connolly, Hillshire Brands)

[Who is the woman who was named CEO of Gilt Groupe in Septemb...](#)

askville.amazon.com › Miscellaneous › Popular News

Askville Question: Who is the woman **who was named CEO of** Gilt Groupe in September? : Popular News.

CEO(woman, Gilt Groupe)

[Tom McKillop - Wikipedia, the free encyclopedia](#)

en.wikipedia.org/wiki/Tom_McKillop

Sir Thomas Fulton Wilson "Tom" McKillop, FRS (born 19 March 1943) is a Scottish chemist, **who was named CEO of** AstraZeneca PLC in 1999 (retired 1 January ...

CEO(scottish chemist, AstraZeneca)

[Harrison adjusts to view from top at First Hawaiian - Pacific Business ...](#)

www.bizjournals.com/.../harrison-adjusts-to-view-from-top-at.html?...

Jan 27, 2012 – Bob Harrison, **who was named CEO of** First Hawaiian Bank on Jan. 1, says he'll spend a lot of time focusing on his people and community ...

CEO(Bob Harrison, First Hawaiian Bank)

Predictive Analysis of Text

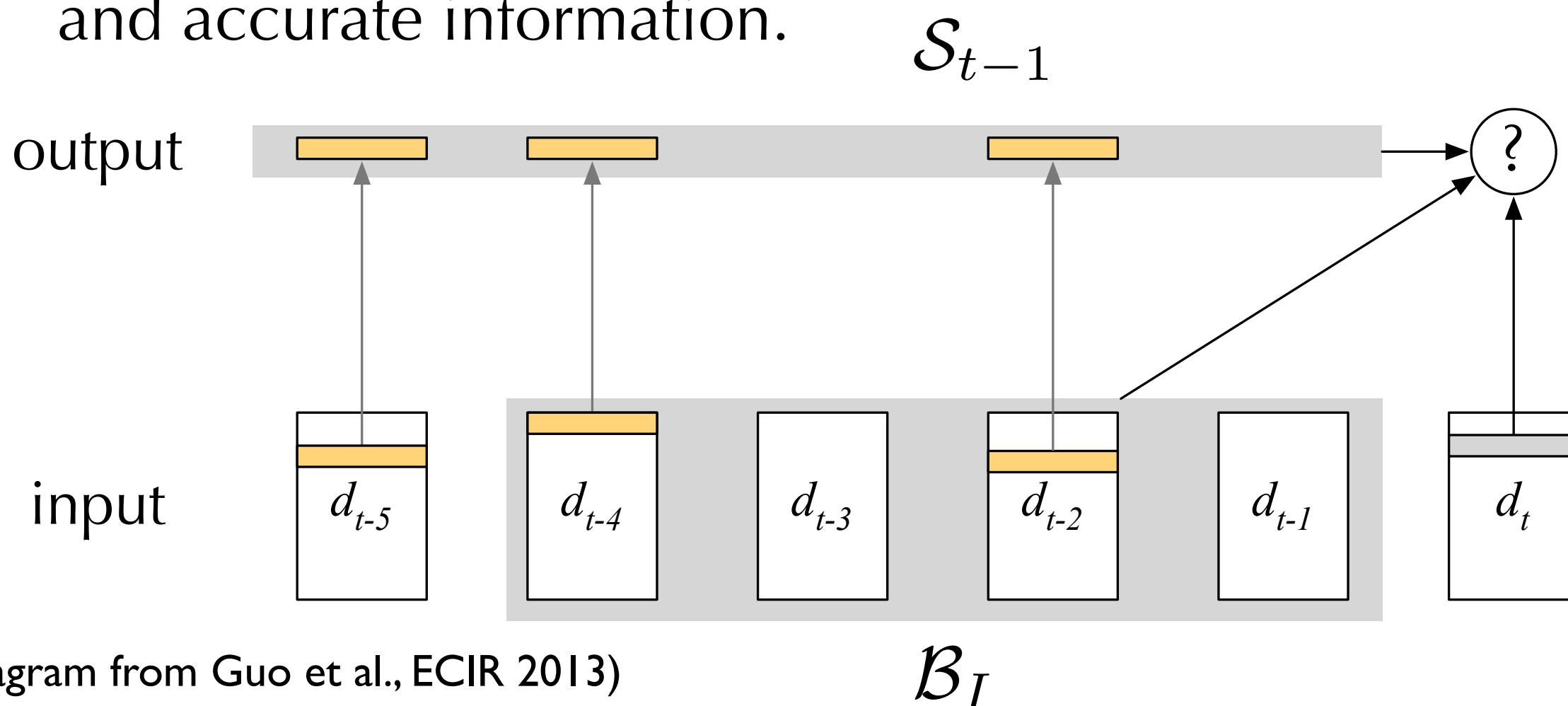
example applications

- **Text-based Forecasting:** monitoring incoming text (e.g., tweets) and making predictions about external, real-world events or trends, for example:
 - ▶ a presidential candidate's poll rating
 - ▶ a company's stock value change
 - ▶ a movie's box office earnings
 - ▶ side-effects for a particular drug
 - ▶ ...

Predictive Analysis of Text

example applications

- **Temporal Summarization:** monitoring incoming text (e.g., tweets) about a news event and predicting whether a sentence should be included in an on-going summary of the event
- Updates to the summary should contain relevant, novel, and accurate information.



(Diagram from Guo et al., ECIR 2013)

Predictive Analysis of Text

example applications

- Detecting other interesting properties of text: [insert your crazy idea here], for example, detecting humorous text:
 - ▶ “Beauty is in the eye of the beholder” not funny
 - ▶ “Beauty is in the eye of the beer holder” funny

(example from Mihalcea and Pulman, 2007)

Course Overview

Road Map

first half of the semester

- Predictive Analysis of Text
 - ▶ Supervised machine learning principles
 - ▶ Text representation
 - ▶ Basic machine learning algorithms
 - ▶ Tools for predictive analysis of text
 - ▶ Experimentation and evaluation
 - ▶ Labeling data and learning models with noisy labels
 - ▶ Feature selection
- Exploratory Analysis of Text
 - ▶ Clustering
 - ▶ Outlier detection (tentative)
 - ▶ Co-occurrence statistics

Road Map

second half of the semester

- Applications
 - ▶ Text classification
 - ▶ Opinion mining
 - ▶ Sentiment analysis
 - ▶ Bias detection
 - ▶ Information extraction
 - ▶ Relation learning
 - ▶ Text-based forecasting
 - ▶ Temporal Summarization
- Is there anything that you would like to learn more about?

Grading

- 40% homework
 - ▶ 10% each
- 15% midterm
- 35% term project
 - ▶ 5% proposal
 - ▶ 10% presentation
 - ▶ 20% paper
- 10% participation

Grading for Graduate Students

- H: 95-100%
- P: 80-94%
- L: 60-79%
- F: 0-59%

Grading for Undergraduate Students

- A+: 97-100%
- A: 94-96%
- A-: 90-93%
- B+: 87-89%
- B: 84-86%
- B-: 80-83%
- C+: 77-79%
- C: 74-76%
- C-: 70-73%
- D+: 67-69%
- D: 64-66%
- D-: 60-63%
- F: $\leq 59\%$

General Outline of Homework

- Given a dataset (i.e., a training and test set), run experiments where you try to predict the target class using different feature representations
- Do error analysis
- Report on what worked, what didn't, and why!
- Answer essay questions about the assignment
 - ▶ These will be associated with the course material

Homework vs. Midterm



- The homework will be more challenging than the midterm. It should be, you have more time.

Term Project

[http://ils.unc.edu/courses/2013_fall/inls613_001/hw/
project_description.pdf](http://ils.unc.edu/courses/2013_fall/inls613_001/hw/project_description.pdf)

Course Tips

- Work hard
- Do the assigned readings
- Do other readings
- Be patient and have reasonable expectations
 - ▶ you're not supposed to understand everything we cover in class during class
- Seek help sooner rather than later
 - ▶ office hours: manning 305, T. and Th. 9:30-10:30 am
 - ▶ questions via email
- Remember the golden rule: no pain, no gain

Questions?