

CSE 590
DATA SCIENCE FUNDAMENTALS

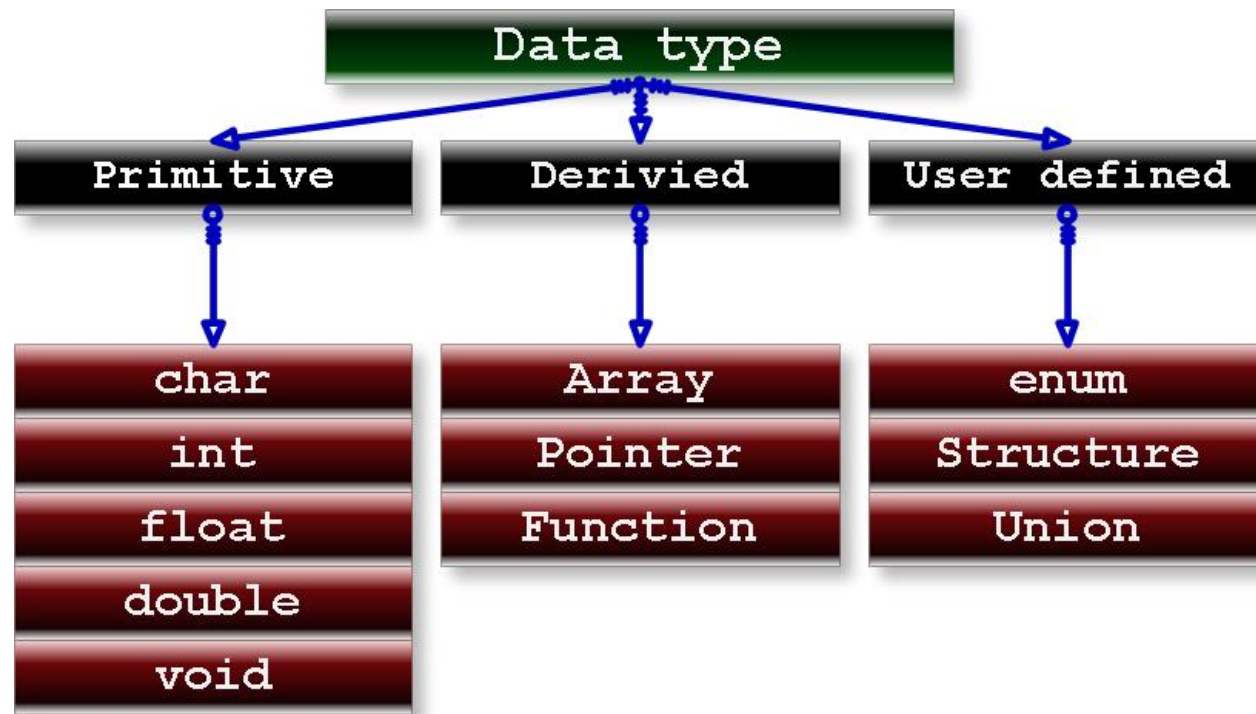
DATA TYPES

KLAUS MUELLER

COMPUTER SCIENCE DEPARTMENT
STONY BROOK UNIVERSITY AND SUNY KOREA

Lecture	Topic	Projects
1	Intro, schedule, and logistics	
2	Data Science components and tasks	
3	Data types	Project #1 out
4	Introduction to R, statistics foundations	
5	Introduction to D3, visual analytics	
6	Data preparation and reduction	
7	Data preparation and reduction	Project #1 due
8	Similarity and distances	Project #2 out
9	Similarity and distances	
10	Cluster analysis	
11	Cluster analysis	
12	Pattern mining	Project #2 due
13	Pattern mining	
14	Outlier analysis	
15	Outlier analysis	Final Project proposal due
16	Classifiers	
17	Midterm	
18	Classifiers	
19	Optimization and model fitting	
20	Optimization and model fitting	
21	Causal modeling	
22	Streaming data	Final Project preliminary report due
23	Text data	
24	Time series data	
25	Graph data	
26	Scalability and data engineering	
27	Data journalism	
	Final project presentation	Final Project slides and final report due

DATA TYPES EVERY CS PERSON KNOWS



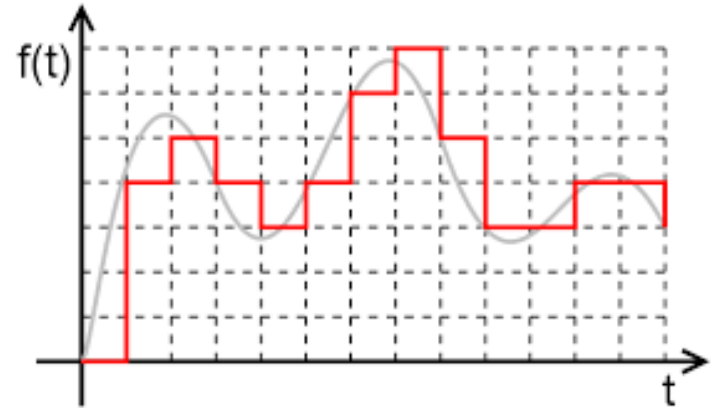
DATA TYPES IN DATA SCIENCE

Source data type
Numeric
Categorical
Text
Time series
Time series
Discrete sequence
Spatial
Graphs
Any type

VARIABLES IN STATISTICS

Numeric variables

- measure a **quantity** as a number
- like: 'how many' or 'how much'
- can be continuous (grey curve)
- or discrete (red steps)



Categorical variables

- describe a **quality** or characteristic
- like: 'what type' or 'which category'
- can be ordinal = ordered, ranked (distances need not be equal)
 - clothing size, academic grades, levels of agreement
- or nominal = not organized into a logical sequence
 - gender, business type, eye color, brand

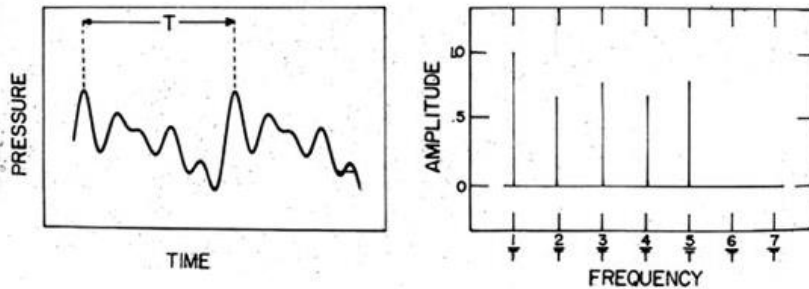
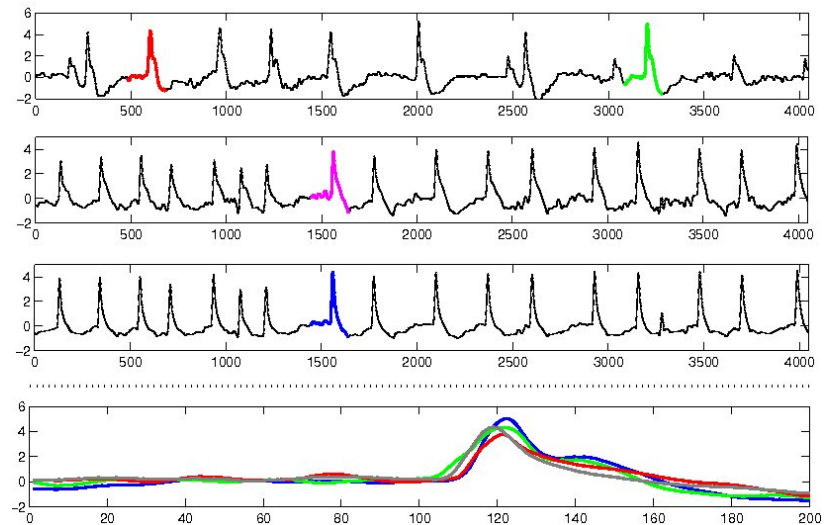
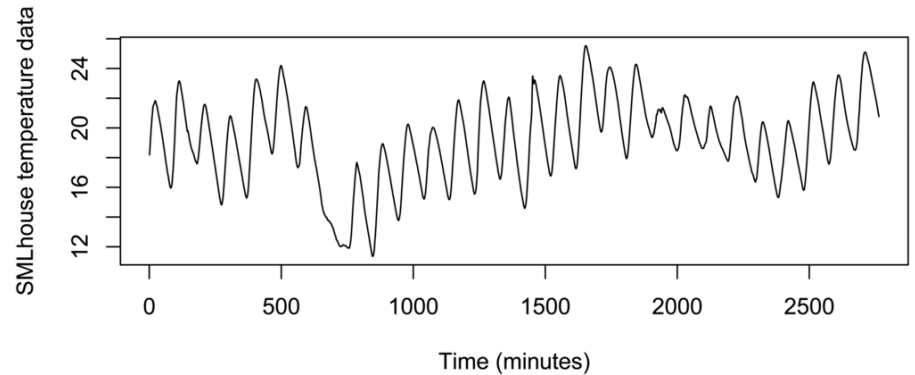
SENSOR DATA

Characteristics

- often large scale
- time series

Feature Analysis

- Fourier transform (FT, FFT)
- Wavelet transform (WT, FWT)
- Motif discovery



Fourier transform

Motif discovery

IMAGE DATA

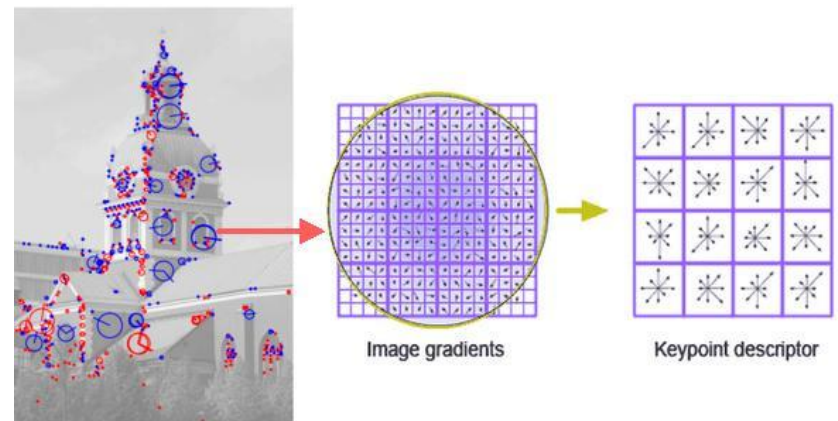
Characteristics

- array of pixels

histograms

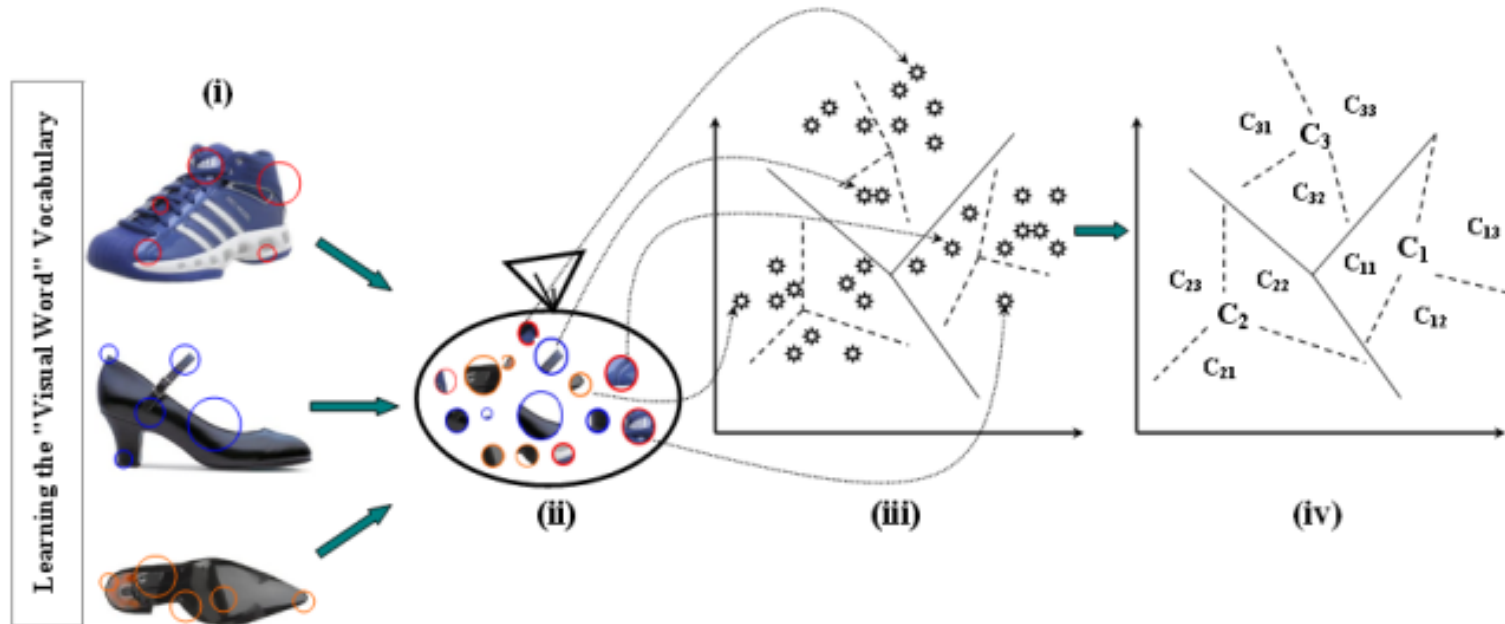
Feature Analysis

- histograms
 - values
 - gradients
- FFT, FWT
- Scale Invariant Feature Transform (SIFT)
- Bag of Features (BoF)
- visual words



SIFT

BAG OF FEATURES (BoF)



BAG OF FEATURES (BoF)

1. Obtain the set of bags of features

- (i) Select a large set of images
- (ii) Extract the SIFT feature points of all the images in the set and obtain the SIFT descriptor for each feature point extracted from each image
- (iii) Cluster the set of feature descriptors for the amount of bags we defined and train the bags with clustered feature descriptors
- (iv) Obtain the visual vocabulary

2. Obtain the BoF descriptor for a given image/video frame

- (v) Extract SIFT feature points of the given image
- (vi) Obtain SIFT descriptor for each feature point
- (vii) Match the feature descriptors with the vocabulary we created in the first step
- (viii) Build the histogram

[More information](#)

VIDEO DATA

Characteristics

- essentially a time series of images

Feature Analysis

- many of the above techniques apply albeit extension is non-trivial



OTHER DATA

Weblogs

- typically represented as text strings in a pre-specified format
- this makes it easy to convert them into multidimensional representation of categorical and numeric attributes

Network traffic

- characteristics of the network packets are used to analyze intrusions or other interesting activity
- a variety of features may be extracted from these packets
 - the number of bytes transferred
 - the network protocol used
 - IP ports used

TEXT DATA

Characteristics

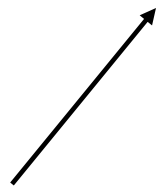
- often raw and unstructured

Feature analysis

- first step is to remove stop words and stem the data
- perform **named-entity recognition** to gain atomic elements
 - identify names, locations, actions, numeric quantities, relations
 - understand the structure of the sentence and complex events
- example:
 - Jim bought 300 shares of Acme Corp. in 2006.
 - [Jim]_{Person} bought [300 shares]_{Quantity} of [Acme Corp.]_{Organiz.} in [2006]_{Time}
- distinguish between
 - application of grammar rules (old style, need experienced linguists)
 - statistical models (Google etc., need big data to build)

DATA TYPES IN DATA SCIENCE

Source data type	Destination data type	Methods
Numeric	Categorical	Discretization
Categorical	Numeric	Binarization
Text	Numeric	Latent semantic analysis (<i>LSA</i>)
Time series	Discrete sequence	<i>SAX</i>
Time series	Numeric multidimensional	<i>DWT, DFT</i>
Discrete sequence	Numeric multidimensional	<i>DWT, DFT</i>
Spatial	Numeric multidimensional	2-d <i>DWT</i>
Graphs	Numeric multidimensional	<i>MDS</i> , spectral
Any type	Graphs	Similarity graph (Restricted applicability)



which you may need to convert into



using these methods

NUMERIC TO CATEGORICAL DATA: DISCRETIZATION

Goal

- divide the ranges of the numeric attribute into φ ranges
- examples:
 - age: [0, 10], [11, 20], [21, 30], ... \rightarrow 0, 1, 2, 3, ...
 - salary: [40,000, 80,000], ... [1,040,000, 1,080,000] \rightarrow 1, ..., 26,...
- what is lost here?
 - distribution information within group
 - data might be distributed unevenly (less for age, more for salary)
- what can we do?
 - store a value that characterizes the distribution
 - use log scaling (when the attribute shows an exponential distribution)
 - use equi-depth range – store the same number of records per bin
 - any disadvantages?

CATEGORICAL TO NUMERIC DATA: BINARIZATION

If a categorical attribute has φ different values, then φ different binary attributes are created

This can have disadvantages

- ordinal data are not spaced by actual distance, whatever the metric
- nominal data have no order and spacing at all
- there are techniques based on correlation (stay tuned)

TEXT TO NUMERIC DATA

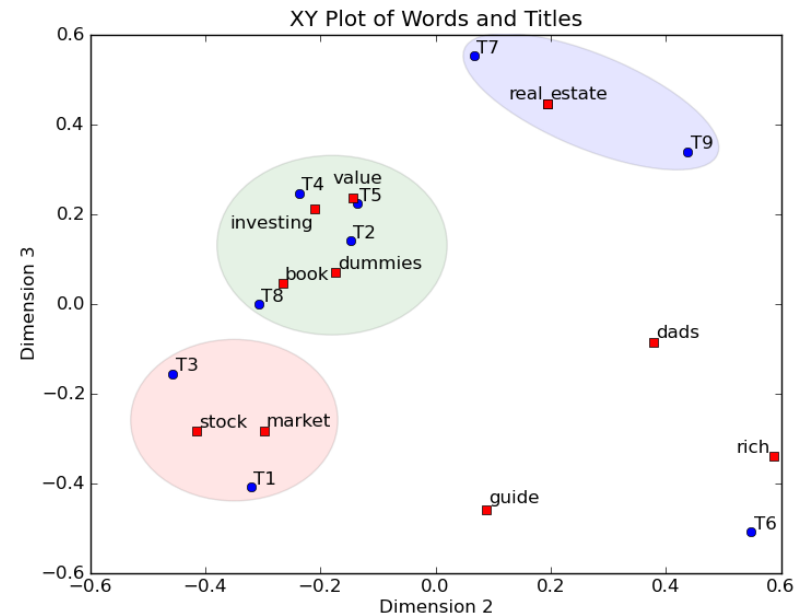
Use Latent Semantic Analysis (LSA)

- also known as Latent Semantic Indexing (LSI)
- turns text into a high-dimensional vector which can be compared

Index Words	Titles								
	T1	T2	T3	T4	T5	T6	T7	T8	T9
book			1	1					
dads						1			1
dummies		1						1	
estate							1		1
guide	1					1			
investing	1	1	1	1	1	1	1	1	1
market	1		1						
real							1		1
rich						2			1
stock	1		1					1	
value				1	1				

Count Matrix

LSA
→



Word/document cluster

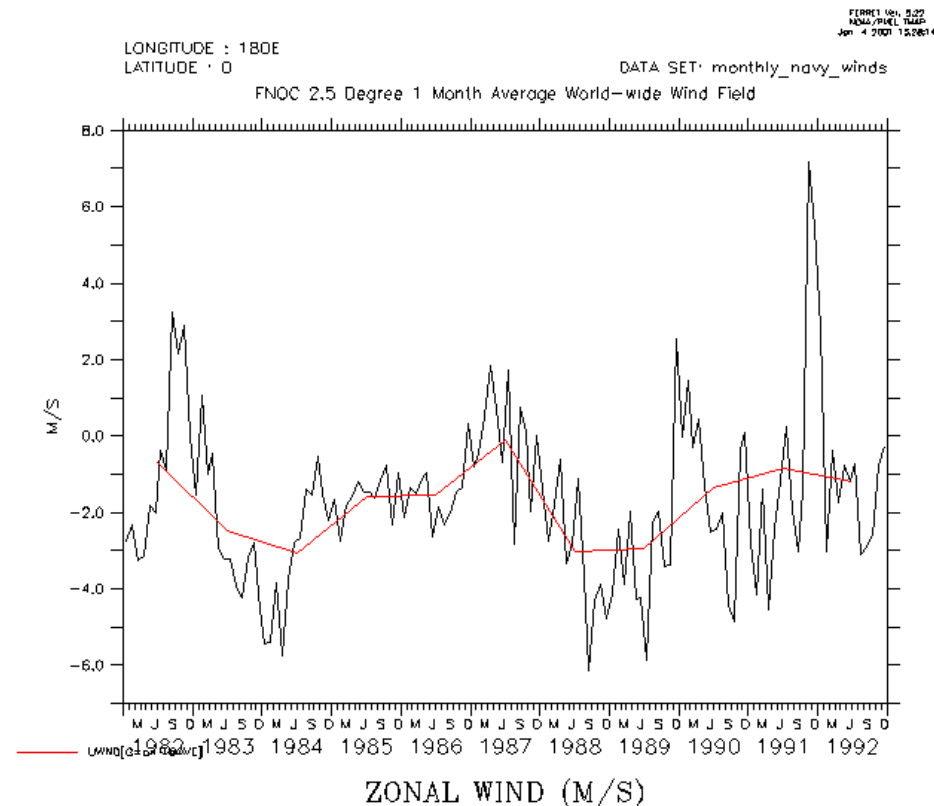
SEQUENCE DATA

Time series to discrete sequence data

- window based averaging
 - eliminates noise in the data
- value-based discretization
 - collapse series into a smaller number of intervals
 - usually done after averaging

Time series to numeric data

- FFT, DFT
- also works for discrete series and spatial data



PROJECT #1

Dataset of NYC bike sharing service

- ~1M records each for July 2013, 2014, 2015
- will link from the course website (under *Labs*)

Attributes:

- Trip Duration
- Start Time & Date, Stop Time & Date
- Start/End Station Name
- Station ID
- Station Latitude/Longitude
- Bike ID
- User Type (24-hour pass, 7-day pass, Annual)
- User Gender, Year of Birth

PROJECT #1 (CONT'D)

Tasks:

- make interesting observations about the data
- keep in mind the nature of the data
- what might be important to the city, provider, users, investors?

Strategies

- use R (next lecture) to analyze the data
- R has visualization tools to present the findings
- can also use D3 (lecture after next)

Deliverables

- comprehensive report with findings (numbers, plots, text narration)
- submit to instructor by email by 9/22