

CSE 332

INTRODUCTION TO VISUALIZATION

VISUAL ANALYTICS BASIC TASKS &
SOLUTION DESIGN

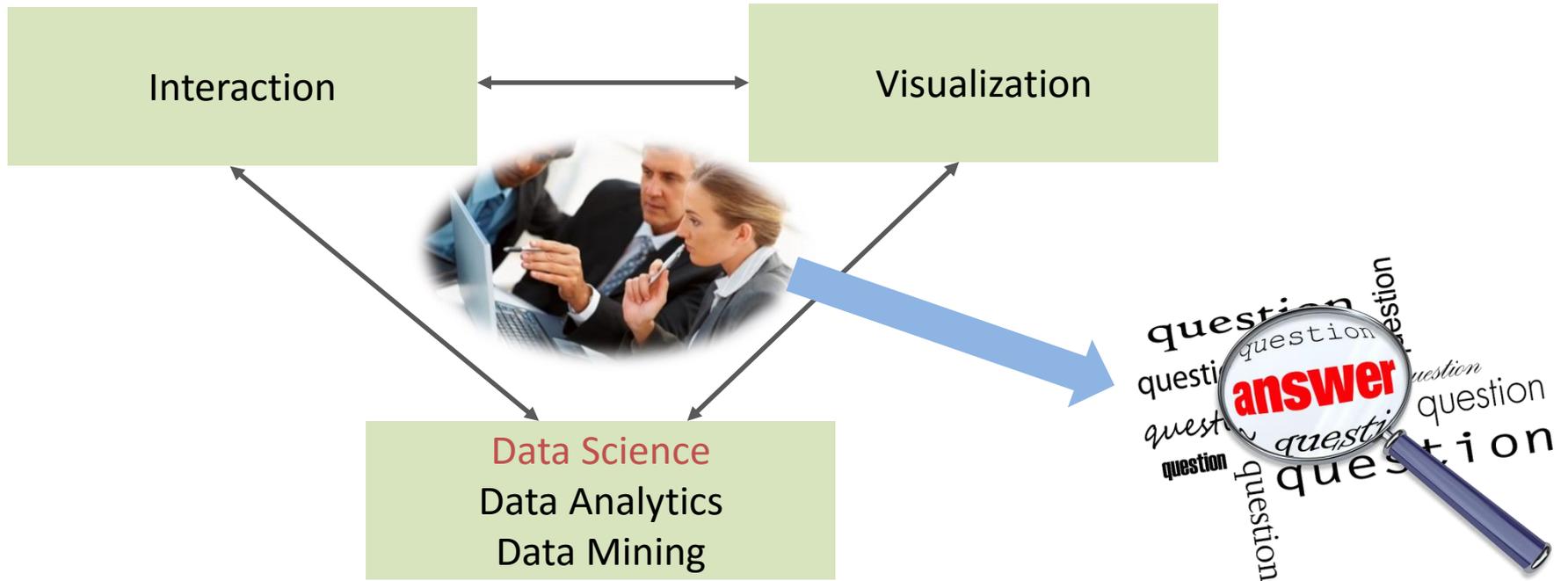
KLAUS MUELLER

COMPUTER SCIENCE DEPARTMENT
STONY BROOK UNIVERSITY

| Lecture | Topic | Projects |
|---------|---|---------------|
| 1 | Intro, schedule, and logistics | |
| 2 | Applications of visual analytics, data types | |
| 3 | Data sources and preparation | Project 1 out |
| 4 | Data reduction, similarity & distance, data augmentation | |
| 5 | Dimension reduction | |
| 6 | Introduction to D3 | |
| 7 | Visual communication using infographics | |
| 8 | Visual perception and cognition | Project 2 out |
| 9 | Visual design and aesthetic | |
| 10 | Cluster analysis | |
| 11 | Visual analytics tasks and design | |
| 12 | High-dimensional data VIS, dimensionality reduction | |
| 13 | Visualization of spatial data: volume visualization intro | Project 3 out |
| 14 | Introduction to GPU programming | |
| 15 | Visualization of spatial data: raycasting, transfer functions | |
| 16 | Illumination and isosurface rendering | |
| 17 | Midterm | |
| 18 | Scientific visualization | |
| 19 | Non-photorealistic and illustrative rendering | Project 4 out |
| 20 | Midterm discussion | |
| 21 | Principles of interaction | |
| 22 | Visual analytics and the visual sense making process | |
| 23 | Visualization of graphs and hierarchies | |
| 24 | Visualization of time-varying and streaming data | Project 5 out |
| 25 | Maps | |
| 26 | Memorable visualizations, visual embellishments | |
| 27 | Evaluation and user studies | |
| 28 | Narrative visualization, storytelling, data journalism, XAI | |

PROLOGUE

Overall definition of visual analytics



- What are the fundamental tasks of data science?
- How can humans assist in these?
- How can humans benefit from these?

FUNDAMENTAL TASKS IN DATA SCIENCE

TASK #1: CLASSIFICATION

Predict which class a member of a certain population belongs to

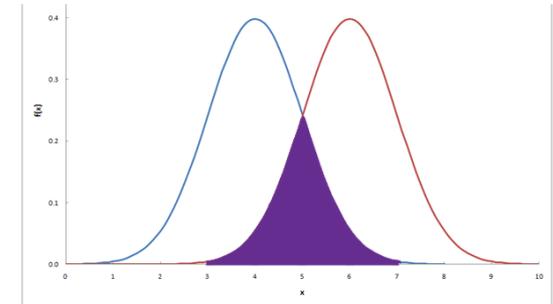
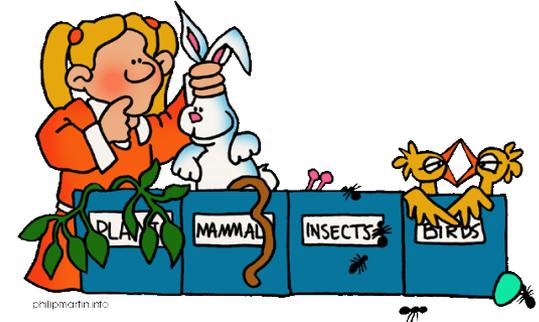
- absolute
- probabilistic

Require a classification model

- absolute
- probabilistic (likelihood)

Scoring with a model

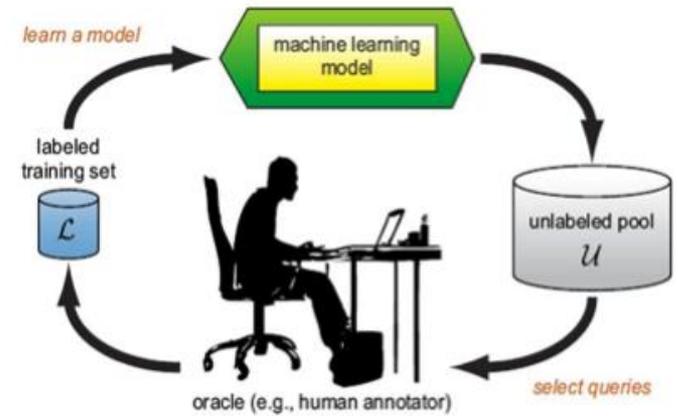
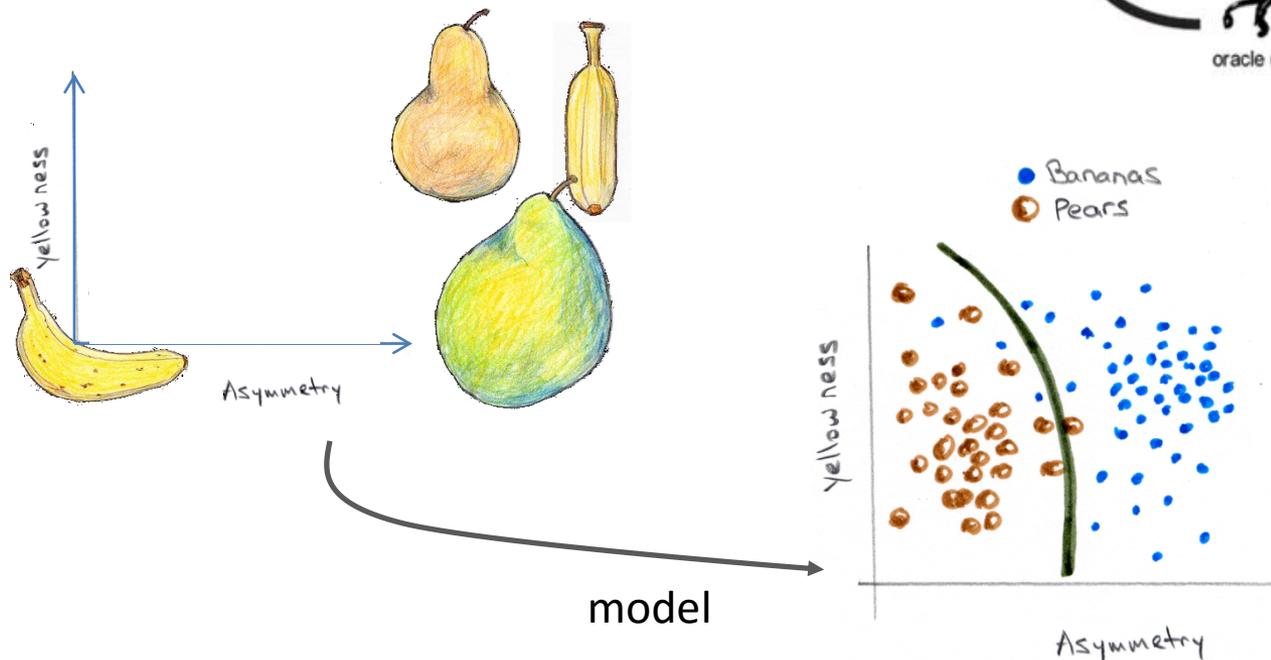
- each population member gets a score for a particular class/category
- sort each class or member scores to assign
- scoring and classification are related



CLASSIFICATION: THE HUMAN FACTOR

Supervised learning

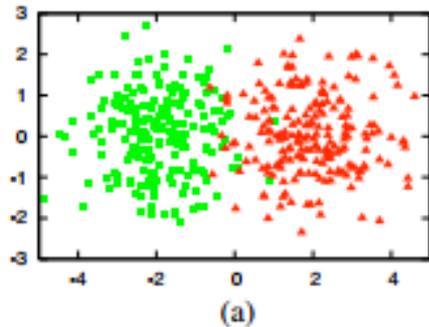
- human labels the samples
- find a good feature vector
- build the classification model



CLASSIFICATION: THE HUMAN FACTOR

Active supervised learning

- human labels the samples
- but while samples are often abundant, labeling can be expensive
- active learning → only label the samples critical to the model

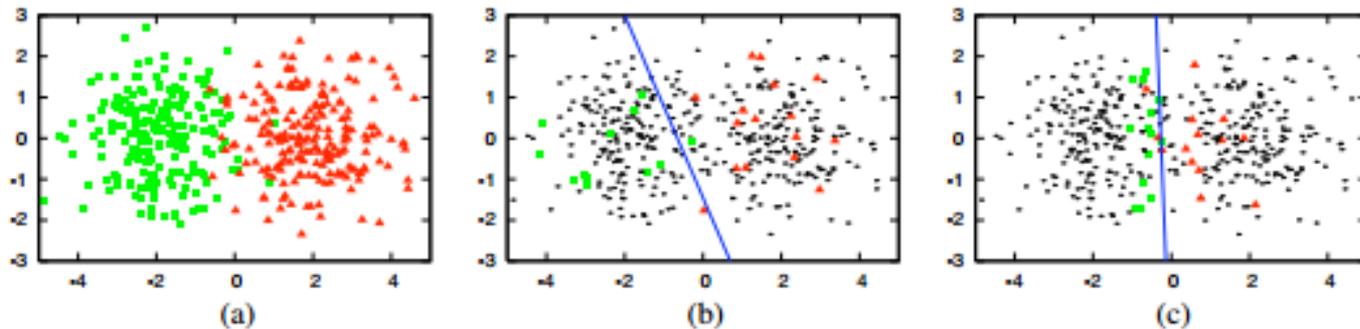


- Assume a toy data set of 400 instances, evenly sampled from two class Gaussians, visualized in 2D feature space.
- Learn a logistic regression model by training it with 30 labeled instances randomly drawn from the problem domain (70% accuracy)
- Learn a logistic regression model by training it with 30 actively queried instances using uncertainty sampling (90%)

CLASSIFICATION: THE HUMAN FACTOR

Active supervised learning

- human labels the samples
- but while samples are often abundant, labeling can be expensive
- active learning → only label the samples critical to the model

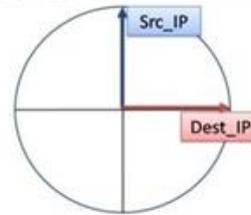
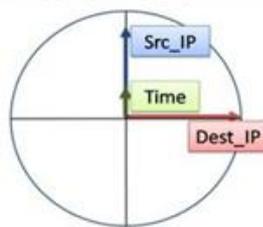
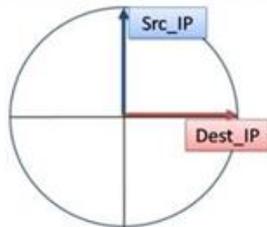
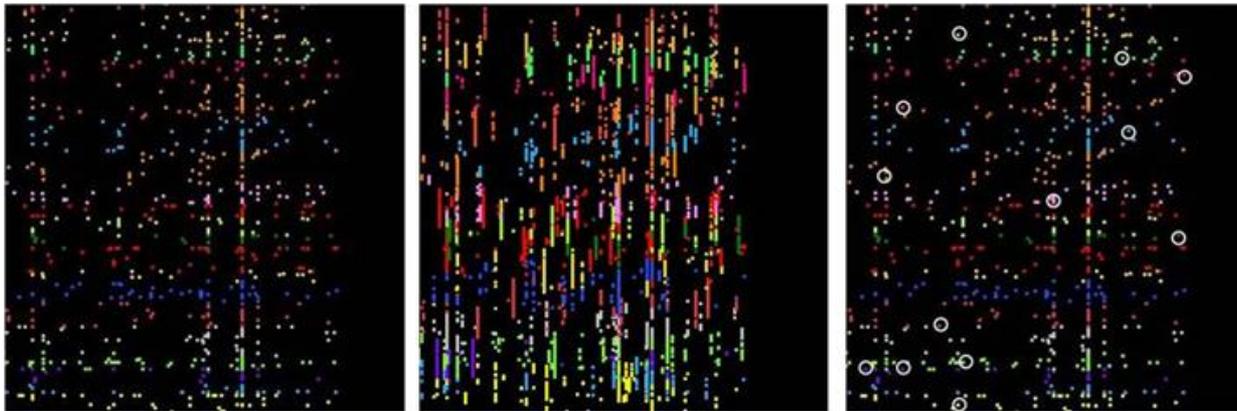


- Assume a toy data set of 400 instances, evenly sampled from two class Gaussians, visualized in 2D feature space.
- Learn a logistic regression model by training it with 30 labeled instances randomly drawn from the problem domain (70% accuracy)
- Learn a logistic regression model by training it with 30 actively queried instances using uncertainty sampling (90%)

APPLICATION: VISUAL MODEL LEARNING

Simple example: network traffic analysis

- the (very large) data set consists of a 1-hour snapshot of internet packets
- goal is to learn the concept 'webpage load'



Mark good
examples

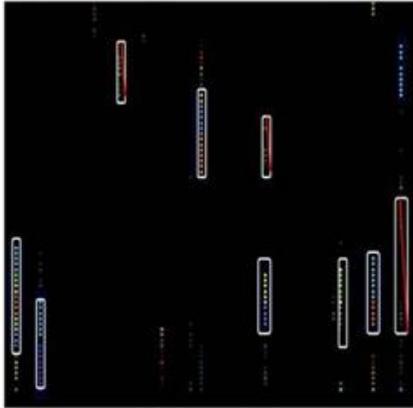
VISUAL MODEL LEARNING: SET INITIAL RULE

Use Inductive Logic Programming (Prolog) to formulate initial model (rule):

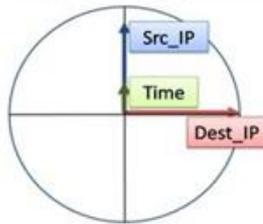
```
webpage_load(X) :-  
    same_src_ips(X), same_dest_ips(X), same_src_port(X, 80)
```

VISUAL MODEL LEARNING: VERIFY INITIAL RULE

Now we classify other data points with this rule and visualize



Mark negative
examples



VISUAL MODEL LEARNING: REFINE INITIAL RULE

Marking negative examples yields updated/refined rule:

```
webpage_load(X) :-  
same_src_ips(X), same_dest_ips(X), same_src_port(X, 80),  
timeframe_upper(X, 10), length(X, L), greaterthan(L, 8).
```

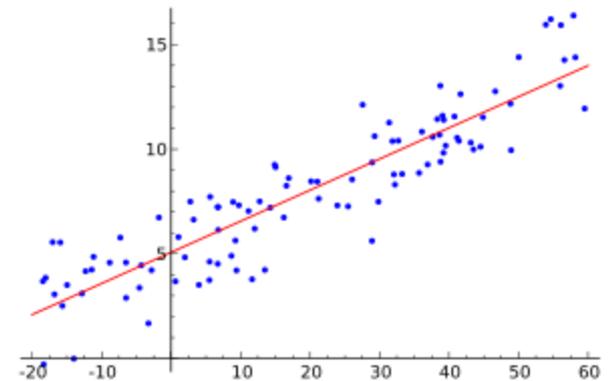
here: must contain at least 8 packets and be within a time frame of 10

TASK #2: REGRESSION

Regression = value estimation

Fit the data to a function

- often linear, but does not have to be
- quality of fit is decisive



Regression vs. classification

- classification predicts that something will happen
- regression predicts how much of it will happen

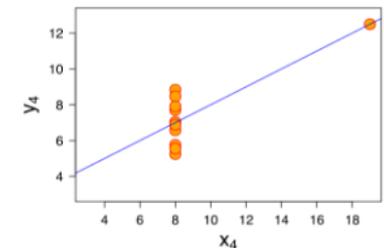
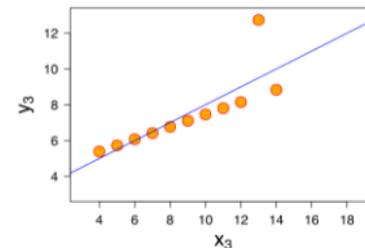
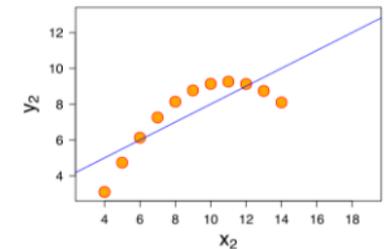
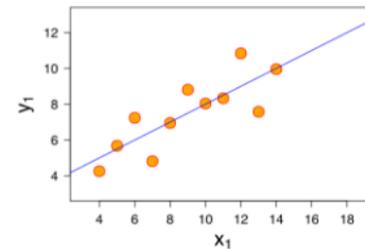
Human factor:

- identify possible outliers

ANSCOMBE QUARTET

Visualization of statistics results is important

| I | | II | | III | | IV | |
|----|-------|----|------|-----|-------|----|------|
| x | y | x | y | x | y | x | y |
| 10 | 8.04 | 10 | 9.14 | 10 | 7.46 | 8 | 6.58 |
| 8 | 6.95 | 8 | 8.14 | 8 | 6.77 | 8 | 5.76 |
| 13 | 7.58 | 13 | 8.74 | 13 | 12.74 | 8 | 7.71 |
| 9 | 8.81 | 9 | 8.77 | 9 | 7.11 | 8 | 8.84 |
| 11 | 8.33 | 11 | 9.26 | 11 | 7.81 | 8 | 8.47 |
| 14 | 9.96 | 14 | 8.10 | 14 | 8.84 | 8 | 7.04 |
| 6 | 7.24 | 6 | 6.13 | 6 | 6.08 | 8 | 5.25 |
| 4 | 4.26 | 4 | 3.10 | 4 | 5.39 | 19 | 12.5 |
| 12 | 10.84 | 12 | 9.13 | 12 | 8.15 | 8 | 5.56 |
| 7 | 4.82 | 7 | 7.26 | 7 | 6.42 | 8 | 7.91 |
| 5 | 5.68 | 5 | 4.74 | 5 | 5.73 | 8 | 6.89 |



| Property | Value |
|--|---|
| Mean of x in each case | 9 (exact) |
| Sample variance of x in each case | 11 (exact) |
| Mean of y in each case | 7.50 (to 2 decimal places) |
| Sample variance of y in each case | 4.122 or 4.127 (to 3 decimal places) |
| Correlation between x and y in each case | 0.816 (to 3 decimal places) |
| Linear regression line in each case | $y = 3.00 + 0.500x$ (to 2 and 3 decimal places, respectively) |

Same statistics
Very different data

TASK #3: SIMILARITY MATCHING

Identify similar individuals based on data known about them

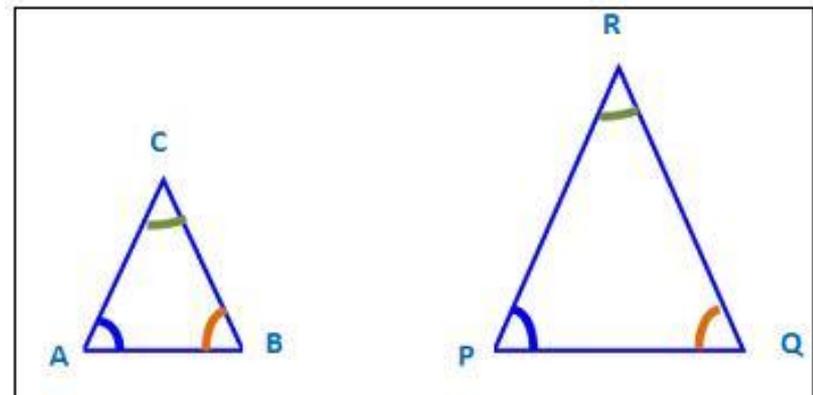
- need a measure of similarity
- features that define similarity
- characteristics

Similarity often part of

- classification
- regression
- clustering

Human factor

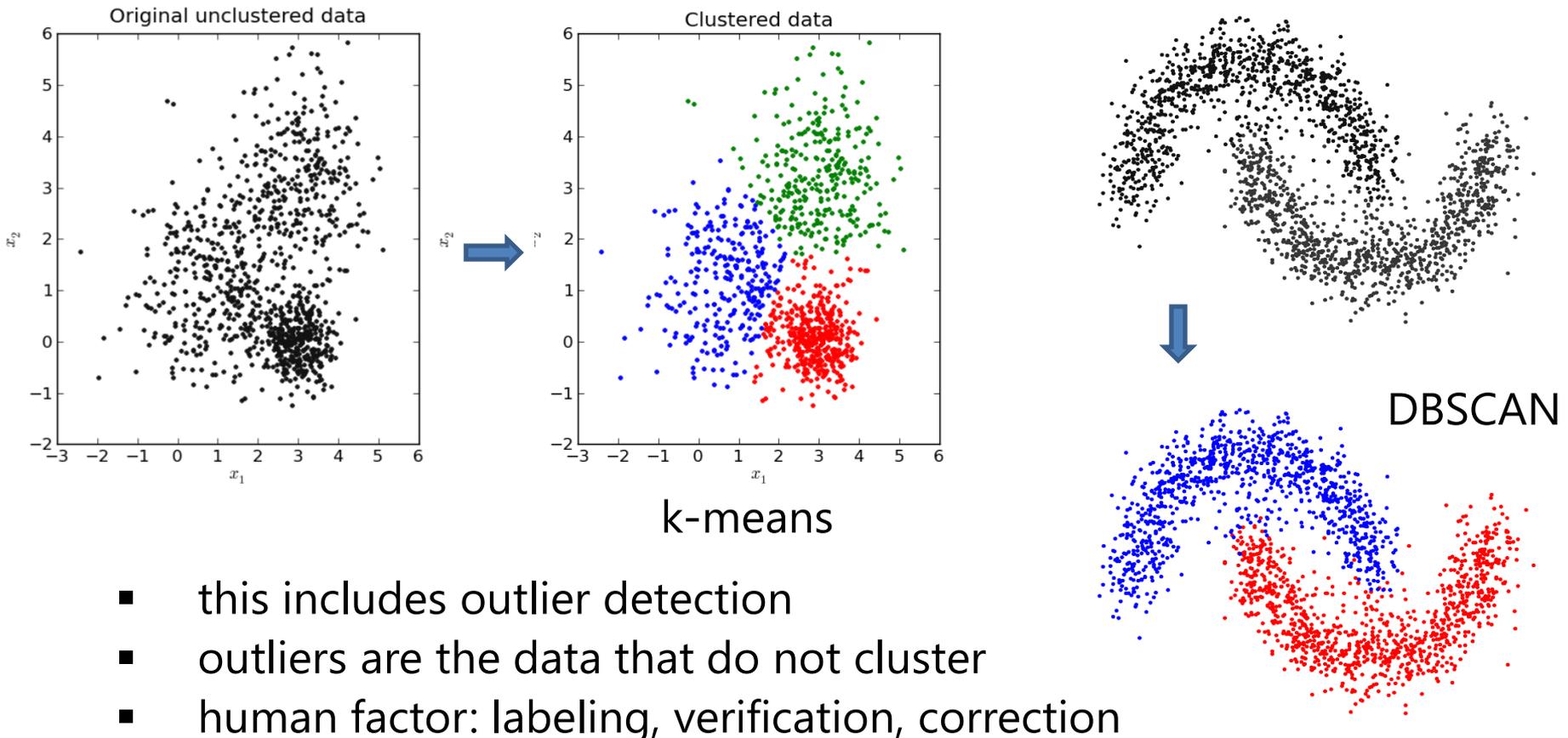
- similar to supervised learning
- identify effective features



TASK #4: CLUSTERING

Group individuals in a population together by their similarity

- preliminary domain exploration to see which natural groups exist



- this includes outlier detection
- outliers are the data that do not cluster
- human factor: labeling, verification, correction

TASK #6: PROFILING

Also known as behavior description

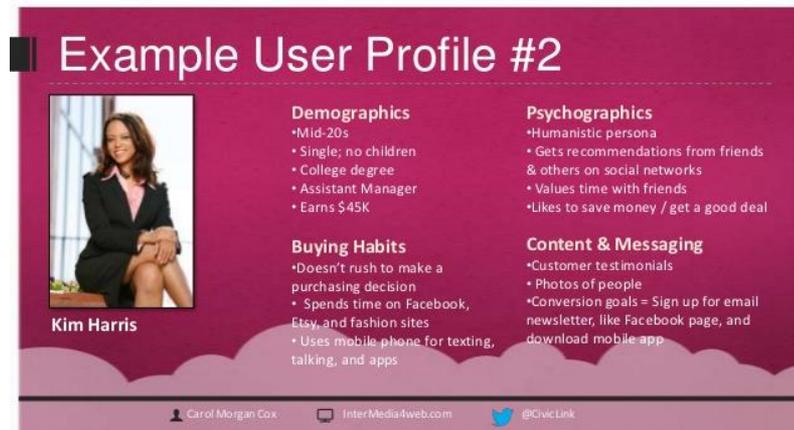
- attempts to **characterize** the typical behavior of an individual, group, or population

Often used to establish behavioral norms for **anomaly detection**

- fraud detection
- intrusion detection

Examples:

- credit card fraud
- airport security



Example User Profile #2

Demographics

- Mid-20s
- Single; no children
- College degree
- Assistant Manager
- Earns \$45K

Psychographics

- Humanistic persona
- Gets recommendations from friends & others on social networks
- Values time with friends
- Likes to save money / get a good deal

Buying Habits

- Doesn't rush to make a purchasing decision
- Spends time on Facebook, Etsy, and fashion sites
- Uses mobile phone for texting, talking, and apps

Content & Messaging

- Customer testimonials
- Photos of people
- Conversion goals = Sign up for email newsletter, like Facebook page, and download mobile app

Kim Harris

Carol Morgan Cox | InterMedia4web.com | @CivicLink

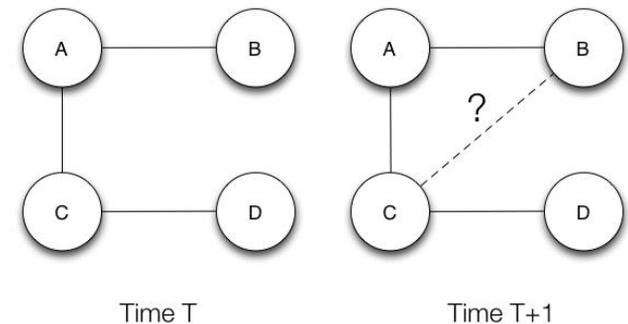
Human factor:

- labeling, verification, correction

TASK #7: LINK PREDICTION

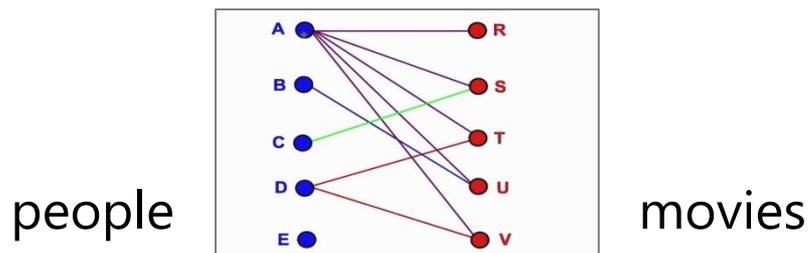
Predict connections between data items

- usually works within a graph
- predict missing links
- estimate link strength



Applications

- in recommendation systems
- friend suggestion in Facebook (social graph)
- link suggestion in LinkedIn (professional graph)
- movie suggestion in Netflix (bipartite graph people – movies)



Human factor:

- labeling
- verification
- correction

TASK #8: DATA REDUCTION

Take a large dataset and substitute it with a smaller one

- keep loss of information minimal
- clustering and cleaning
- importance sampling
- dimension reduction
- data abstraction
- big data → small data
- find *latent* variables



Example for latent variable – Movie *Taste*

- not directly measurable – latent variable
- derive from movie viewing preferences
- can reveal genre, etc.

Human factor:

- labeling
- verification
- correction

TASK #9: CAUSAL MODELING

Understand what events or actions influence others



Different from predictive modeling

- tries **to explain why** the predictive model worked (or not)

Potentially unreliable when done from observational data

- conducting a targeted experiment is better, but often impossible
- have to work with observational (often anecdotal data)
- hence there is a clear human factor: verify the model, correct it, edit it

Builds on counterfactual analysis

- an event is causal if mutating it will lead to undoing the outcome
- “If only I hadn't been speeding, my car wouldn't have been wrecked”
- downward vs. upward counterfactual thinking
- can explain happiness of bronze medalists vs. silver medalists
- just making the grade vs. just missing the grade

CASE STUDY: WHAT CAUSES LOW MPG

THE CAR DATA SET

Consider the salient features of a car (not really big data):

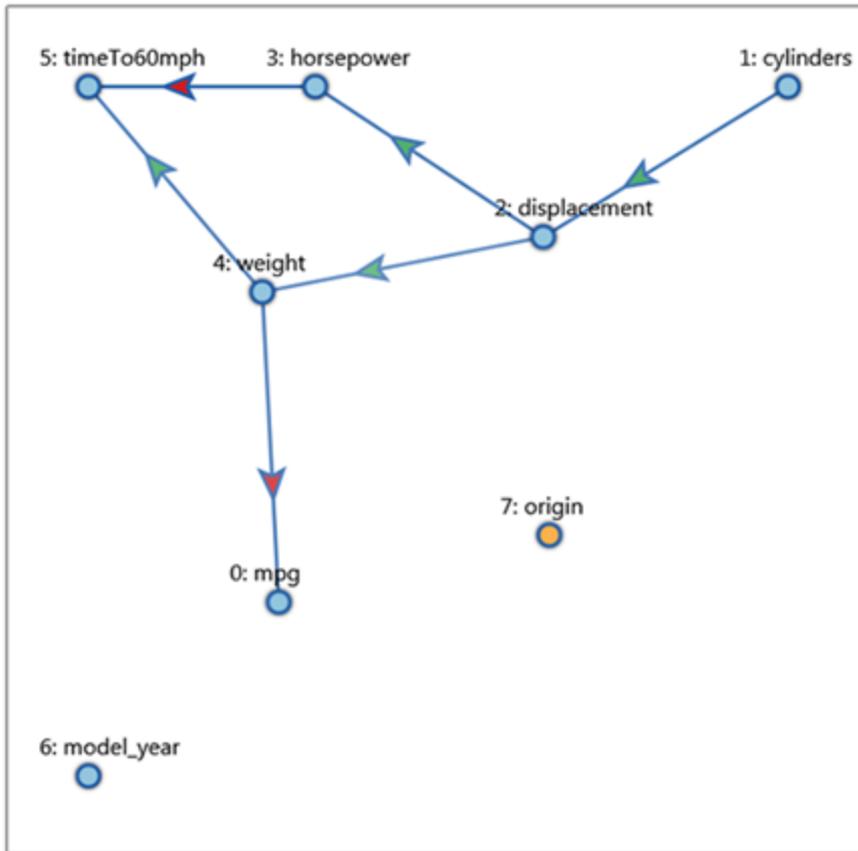
- miles per gallon (MPG)
- top speed
- acceleration (time to 60 mph)
- number of cylinders
- horsepower
- weight
- country origin

400 cars from the 1980s

SHOWN IN A SPREADSHEET

| A1 | | Urban population | | | | | | | | | | | | | | |
|----|------------------------|------------------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|
| | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P |
| 1 | Urban population | 1960 | 1961 | 1962 | 1963 | 1964 | 1965 | 1966 | 1967 | 1968 | 1969 | 1970 | 1971 | 1972 | 1973 | 1974 |
| 2 | Afghanistan | 769308 | 811389 | 855131 | 900646 | 948060 | 997499 | 1053104 | 1110728 | 1170961 | 1234664 | 1302370 | 1391081 | 1483942 | 1579748 | 1676656 |
| 3 | Albania | 494443 | 511637 | 529182 | 547024 | 565117 | 583422 | 601897 | 620508 | 639234 | 658062 | 676985 | 698179 | 719561 | 741149 | 762972 |
| 4 | Algeria | 3293999 | 3513320 | 3737362 | 3969886 | 4216744 | 4483048 | 4644898 | 4822860 | 5015071 | 5218184 | 5429743 | 5618190 | 5813978 | 6017932 | 6231383 |
| 5 | American Samoa | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - |
| 6 | Andorra | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - |
| 7 | Angola | 521205 | 552777 | 585121 | 618345 | 652638 | 688181 | 729595 | 772643 | 817418 | 863993 | 912486 | 982944 | 1056617 | 1133936 | 1215437 |
| 8 | Antigua and Barbuda | 21699 | 21737 | 21878 | 22086 | 22309 | 22513 | 22717 | 22893 | 23053 | 23218 | 23394 | 24046 | 24718 | 25342 | 25826 |
| 9 | Argentina | 15224096 | 15588864 | 15957125 | 16328045 | 16700303 | 17073371 | 17432905 | 17793789 | 18160868 | 18540720 | 18938137 | 19335571 | 19750609 | 20180707 | 20621674 |
| 10 | Armenia | 957974 | 1008899 | 1061551 | 1115546 | 1170414 | 1225785 | 1281346 | 1337060 | 1393199 | 1450241 | 1508526 | 1565054 | 1622558 | 1680709 | 1739019 |
| 11 | Aruba | 24996 | 25514 | 26019 | 26498 | 26941 | 27337 | 27683 | 27984 | 28247 | 28491 | 28726 | 28959 | 29188 | 29409 | 29610 |
| 12 | Australia | 8375329 | 8585577 | 8840666 | 9055650 | 9279777 | 9508980 | 9770529 | 9937118 | 10157212 | 10416192 | 10668471 | 11050785 | 11271606 | 11461308 | 11771589 |
| 13 | Austria | 4560057 | 4589541 | 4621666 | 4653194 | 4685421 | 4715750 | 4754588 | 4778506 | 4798552 | 4817322 | 4849178 | 4871380 | 4904030 | 4932109 | 4939292 |
| 14 | Azerbaijan | 1857673 | 1929429 | 2004258 | 2080816 | 2157307 | 2232355 | 2306310 | 2378380 | 2448728 | 2517815 | 2586000 | 2660687 | 2734631 | 2807879 | 2880491 |
| 15 | Bahamas | 65457 | 69655 | 74179 | 78961 | 83902 | 88918 | 93931 | 98974 | 103944 | 108721 | 113219 | 117339 | 121142 | 124761 | 128393 |
| 16 | Bahrain | 128480 | 133815 | 139791 | 146052 | 152097 | 157596 | 162844 | 167630 | 172373 | 177677 | 183997 | 191379 | 199768 | 209201 | 219678 |
| 17 | Bangladesh | 2761049 | 2947191 | 3141372 | 3344120 | 3556037 | 3777716 | 4047121 | 4329144 | 4624445 | 4933701 | 5257558 | 5710277 | 6184871 | 6682073 | 7202503 |
| 18 | Barbados | 84884 | 85284 | 85761 | 86285 | 86797 | 87259 | 87707 | 88117 | 88526 | 88986 | 89532 | 90518 | 91596 | 92713 | 93796 |
| 19 | Belarus | 2656152 | 2774166 | 2896449 | 3022217 | 3150553 | 3280410 | 3415984 | 3554673 | 3695363 | 3836802 | 3977600 | 4131179 | 4285735 | 4439788 | 4591705 |
| 20 | Belgium | 8435075 | 8489549 | 8548773 | 8620194 | 8709437 | 8796088 | 8865259 | 8924327 | 8968568 | 9003536 | 9040444 | 9086816 | 9134227 | 9175144 | 9217085 |
| 21 | Belize | 49165 | 50608 | 52156 | 53734 | 55226 | 56561 | 57756 | 58820 | 59746 | 60532 | 61186 | 61883 | 62445 | 62984 | 63665 |
| 22 | Benin | 211033 | 229172 | 248065 | 267765 | 288321 | 309788 | 332782 | 366019 | 396065 | 427482 | 460341 | 500355 | 542251 | 586179 | 632320 |
| 23 | Bermuda | 44400 | 45500 | 46600 | 47700 | 48900 | 50100 | 51000 | 52000 | 53000 | 54000 | 55000 | 54600 | 54200 | 53800 | 53400 |
| 24 | Bhutan | 8064 | 8778 | 9526 | 10311 | 11137 | 12010 | 13089 | 14230 | 15445 | 16750 | 18158 | 19926 | 21827 | 23858 | 26008 |
| 25 | Bolivia | 1233398 | 1271250 | 1310294 | 1350615 | 1392328 | 1435536 | 1480255 | 1526529 | 1574517 | 1624419 | 1676370 | 1730434 | 1786553 | 1844596 | 1904355 |
| 26 | Bosnia and Herzegovina | 604204 | 637337 | 671124 | 705395 | 739884 | 774380 | 812856 | 851325 | 890011 | 929301 | 969514 | 1008688 | 1048890 | 1089898 | 1131315 |
| 27 | Botswana | 16240 | 17379 | 18583 | 19855 | 21203 | 22631 | 28191 | 34090 | 40352 | 46995 | 54038 | 61638 | 69689 | 78254 | 87422 |
| 28 | Brazil | 32662018 | 34463344 | 36353068 | 38320171 | 40346703 | 42418482 | 44548227 | 46722996 | 48945984 | 51223962 | 53563179 | 56042505 | 58587770 | 61207586 | 63913385 |
| 29 | Brunei | 35501 | 38753 | 42173 | 45802 | 49699 | 53916 | 58461 | 63355 | 68595 | 74157 | 80024 | 83802 | 87671 | 91616 | 95629 |
| 30 | Bulgaria | 2918659 | 3085061 | 3251675 | 3418610 | 3588246 | 3765058 | 3889518 | 4022040 | 4159890 | 4301340 | 4440270 | 4554810 | 4667059 | 4782931 | 4907107 |
| 31 | Burkina Faso | 221872 | 230199 | 238713 | 247472 | 256558 | 266039 | 275958 | 286311 | 297074 | 308196 | 319642 | 332556 | 345877 | 359655 | 373966 |
| 32 | Burundi | 58810 | 61055 | 63344 | 65696 | 68137 | 70683 | 73370 | 76186 | 79034 | 81779 | 84324 | 90879 | 97308 | 103757 | 110494 |
| 33 | Cambodia | 559631 | 578678 | 598248 | 618631 | 640243 | 663272 | 747219 | 835638 | 927177 | 1019449 | 1110079 | 962037 | 806676 | 645287 | 479631 |
| 34 | Cameroon | 751711 | 801009 | 852578 | 906523 | 962928 | 1021891 | 1088521 | 1158289 | 1231375 | 1307967 | 1388275 | 1522958 | 1664410 | 1813278 | 1970385 |
| 35 | Canada | 12375125 | 12764121 | 13145207 | 13536503 | 13941055 | 14345262 | 14727261 | 15108962 | 15470875 | 15800439 | 16142268 | 16381341 | 16640381 | 16920220 | 17221765 |
| 36 | Cape Verde | 32791 | 34353 | 35972 | 37672 | 39487 | 41435 | 43592 | 45884 | 48200 | 50383 | 52314 | 54103 | 55620 | 56940 | 58184 |
| 37 | Cayman Islands | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - |
| 38 | Central African Rep. | 302157 | 317715 | 333986 | 351001 | 368787 | 387357 | 408129 | 429825 | 452326 | 475441 | 499036 | 526414 | 554452 | 583376 | 613530 |
| 39 | Chad | 198777 | 213406 | 228652 | 244499 | 260903 | 277834 | 305390 | 333898 | 363523 | 394530 | 427153 | 467662 | 510348 | 554973 | 601045 |
| 40 | Channel Islands | 42565 | 42665 | 42792 | 42941 | 43102 | 43269 | 43437 | 43604 | 43765 | 43916 | 44051 | 44208 | 43987 | 43907 | 43762 |

SEEKING THE CAUSE OF LOW MPG



Isolating MPG

The Visual Causality Analyst

Choose Dataset

Auto MPG.rds

Selected Variables:

MPG Cylinders
Displacement Horsepower
Weight TimeTo60MPH
ModelYear Origin

Significant Level

0.1 0.05 0.01

Show Node ID

Parameterized

Data Scaling Method

none

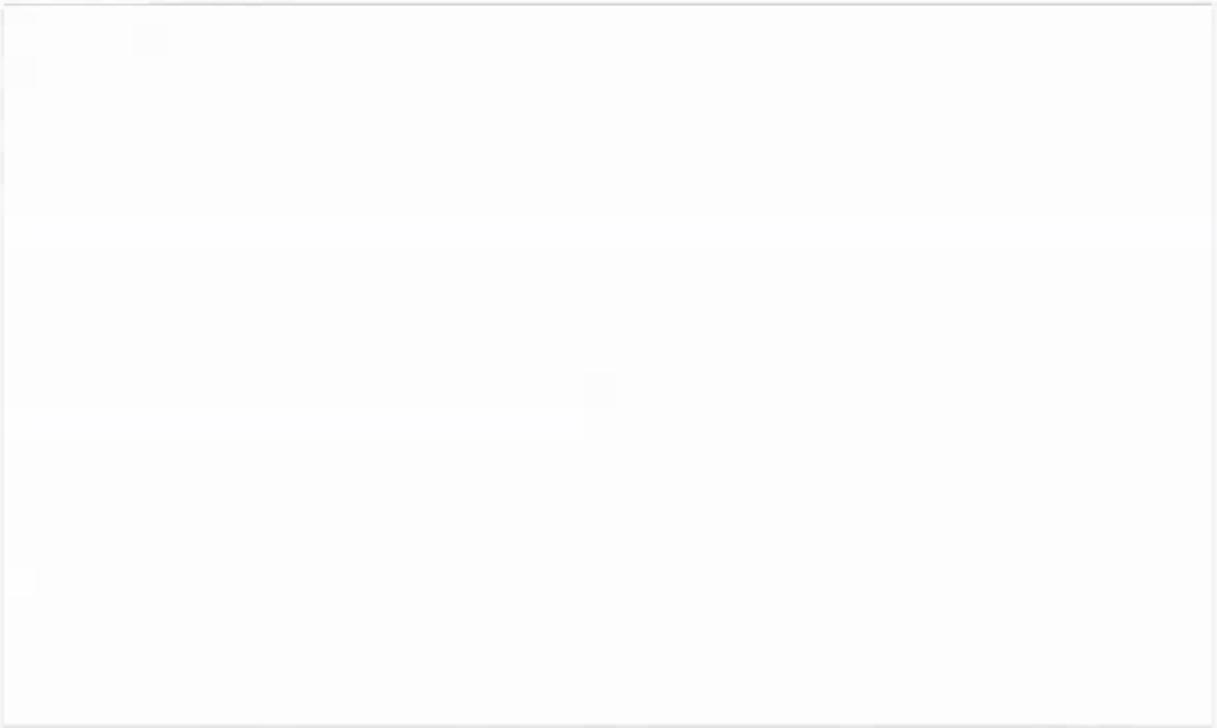
standardize

normalize

Alternative Models

> Infer Causal Model

Causality Viz Data Bracketing



Source: MPG

Target: MPG

Create
Direct

Reverse
Remove

Coefficient Threshold
0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1

Download Graph

[Graph Model Info.]

[Clicked Vertex Info.]

[Clicked Edge Info.]

How To DESIGN A VISUAL ANALYTICS SOLUTION

Use the nested model

- devised by Tamara Munzner (UBC)
- M. Meyer, M. Sedlmair, P. Quinan, T Munzner, "The nested blocks and guidelines model," *Information Visualization*, 2013

STEP 1: CHARACTERIZE THE PROBLEM

Define the tasks, data, workflow of target users

- the tasks are usually described in domain terms
- finding and eliciting the requirements is notoriously hard
- observe how domain users work and perform their tasks
- observe the pains they are having
- what are the limitations?
- what is currently impossible, slow, or tedious?

domain problem characterization

STEP 2: ABSTRACT INTO A DESIGN

Map from domain vocabulary/concerns to abstraction

- may require some sort of transformation
- data and types are described in abstract terms
- numeric tables, relational/network, spatial, ...
- tasks and operations described in abstract terms
- generic activities: sort, filter, correlate, find trends/outliers...

domain problem characterization

data/operation abstraction design

STEP 2: ENCODE INTO A VISUALIZATION

Visual encoding

- how to best show the data (also pay tribute to aesthetics)
- bar/pie/line charts, parallel coordinates, MDS plot, scatterplot, tree map, network, etc.

Interaction design

- how to best support the intent a user may have
- select, navigate, order, brush, ...

domain problem characterization

data/operation abstraction design

encoding/interaction technique design

STEP 4: DESIGN AN ALGORITHM

Well-studied computer science problem

- create efficient algorithms
- should support human interaction
- else it would not comply with key principle of visual analytics

domain problem characterization

data/operation abstraction design

encoding/interaction technique design

algorithm design

APPLICATION EXAMPLE

Let use the causality analyzer framework just presented

- use the car design example

Domain problem characterization

- how to design a faster car without elevating gas consumption

Data/operation abstraction design

- determine how the different car parameters depend on one another
- collect data of different car models and compute a causal network

Encoding/interaction technique design

- draw graph where parameters are nodes and causal links are edges
- provide interactions that allows users to test causal links and compute a score

Algorithm design

- Partial correlation followed by causal inferencing/conditioning
- Bayesian Information Criterion (BIC) to model Occam's Razor

ANOTHER APPLICATION EXAMPLE

How the iPhone came about

- domain problem characterization
- data/operation abstraction design
- encoding/interaction technique design
- algorithm design

June 29, 2007



GAUGE SUCCESS

threat: wrong problem

validate: observe and interview target users

threat: bad data/operation abstraction

threat: ineffective encoding/interaction technique

validate: justify encoding/interaction design

threat: slow algorithm

validate: analyze computational complexity

implement system

validate: measure system time/memory

validate: qualitative/quantitative result image analysis

[test on **any** users, informal usability study]

validate: lab study, measure human time/errors for operation

validate: test on **target** users, collect anecdotal evidence of utility

validate: field study, document human usage of deployed system

validate: observe adoption rates

GAUGE SUCCESS

Validate along the way and refine

- formative user study

Extend to general user studies of the final design

- summative user study
- laboratory study
- smaller number of subjects but can use speak aloud protocol
- crowd-sourced via internet
- potentially greater number of subjects to yield better statistics but can be superficial

We will discuss evaluation studies later