DATA SYSTEMS FOR HUMAN DATA INTERACTION

Remco Chang

Associate Professor Computer Science, Tufts University











"The computer is incredibly fast, accurate, and stupid. Man is unbelievably slow, inaccurate, and brilliant. The marriage of the two is a force beyond calculation."

-Leo Cherne, 1977

(often attributed to Albert Einstein)

WHICH MARRIAGE?



[v]alt

WHICH MARRIAGE?



40M PRACTICE INTO RESERANCE

Visual epresentation

Production, Presentation &

Analytical Reasoning

Data Representations & Transformations

VISUAL ANALYTICS = HUMAN + COMPUTER

• Visual analytics is "the science of analytical reasoning facilitated by visual interactive interfaces."¹

7/53

By definition, it is a collaboration between human and computer to solve data problems.



VISUAL ANALYTICS LAB AT TUFTS (VALT)





Jordan Crouser Prof, Smith

Evan Peck

Prof, Bucknell



Lane Harrison Prof, WPI



Eli Brown Prof, DePaul



Beste Yuksel Prof, USF



Alvitta Ottley Prof, WashU



Leilani Battle



Marianne Procopio **MIT** Lincoln Lab



Shah Humayoun Prof, SF State



Ab Mosca PhD student



Michael Behrisch Prof, Utrecht



Ashley Suh PhD student



Dan Afergan

Google Research

Gabriel Appleby PhD student



Mike Curry **Draper Labs**



Brian Montambault PhD student



Dylan Cashman Novartis Research



Wenbo Tao PhD student (MIT)

VISUAL ANALYTICS LAB AT TUFTS (VALT)

- User Modeling (HCI)
 - Complexity of Human Computation (Crouser)
 - Individual Differences (Ottley, Peck)
 - Perceptual Modeling (Harrison, Yang)
 - Adaptive Systems and BCI (Yuksel, Afergan)
 - Decision Making (Mosca, Ottley)
- Interactive Machine Learning (ML)
 - Learning from User Interactions (Brown)
 - System for Model Selection (Cashman)
 - Projection Methods (Appleby)
 - Explainable Anomaly Detection (Montambault)
- Database Systems (DB)
 - Predictive Prefetching (Battle)
 - Progressive Data Streaming (Procopio)
 - Interactive Pan-Zoom (Tao)



- Financial Fraud Analysis
 - Bank of America
- Political Simulation
 - Agent-based analysis
- Bridge Maintenance
 - Exploring inspection reports
- Biomechanical Motion
 - Interactive motion comparison
- High-D Data Analysis
 - Dis-Function: highdimensional metric learning



R. Chang et al., WireVis: Visualization of Categorical, Time-Varying Data From Financial Transactions. IEEE VAST, 2007

- Financial Fraud AnalysisBank of America
- Political Simulation
 - Agent-based analysis
- Bridge Maintenance
 - Exploring inspection reports
- Biomechanical Motion
 - Interactive motion comparison
- High-D Data Analysis
 - Dis-Function: highdimensional metric learning





R. Chang et al., Two Visualization Tools for Analysis of Agent-Based Simulations in Political Science. IEEE CG&A, 2012

- Financial Fraud Analysis
 - Bank of America
- Political Simulation

 Agent-based analysis
- Bridge Maintenance
 - Exploring inspection reports
- Biomechanical Motion
 - Interactive motion comparison
- High-D Data Analysis
 - Dis-Function: highdimensional metric learning







- Financial Fraud Analysis
 - Bank of America
- Political Simulation
 - Agent-based analysis
- Bridge Maintenance
 - Exploring inspection reports
- Biomechanical Motion
 - Interactive motion comparison
- High-D Data Analysis
 - Dis-Function: highdimensional metric learning







R. Chang et al., Interactive Coordinated Multiple-View Visualization of Biomechanical Motion Data, IEEE Vis (TVCG) 2009.

- Financial Fraud Analysis
 - Bank of America
- Political Simulation
 - Agent-based analysis
- Bridge Maintenance
 - Exploring inspection reports
- Biomechanical Motion
 - Interactive motion comparison
- High-D Data Analysis
 - Dis-Function: highdimensional metric learning



R. Chang et al., Dis-function: Learning Distance Functions Interactively, IEEE VAST 2011.

LESSONS LEARNED

• What do all good marriages have in common?

Collaboration

- Computers and humans need to contribute equally to the task at hand.
- Data systems need to perform complex machine learning tasks with minimal guidance from the user
 - Snowcat
 - "Unprojection"

Communication

- Computers and humans need to "talk" to each other quickly and effortlessly.
- Data systems need to respond to a user's interactions within 500ms
 - ForeCache, Kyrix
 - NeuralCubes





[v]alt



Mike Gleicher (Wisconsin)



John Stasko

(GT)



(GT)

Dylan Cashman (Tufts)



Shah Rukh Humayoun (Tufts, (Wisconsin) SF State)



Florian

Heimerl





Subhajit Das (GT)

(GT)

HUMAN-ML COLLABORATION: VISUAL ANALYTICS WITH AUTOML

R. Chang et al., A User-based Visual Analytics Workflow for Exploratory Model Analysis, EuroVis 2019 R. Chang et al., Defining an Analysis: A Study of Client-Facing Data Scientists, EuroVis 2019 (Short) R. Chang et al., Ablate, Variate, and Contemplate: Visual Analytics for Discovering Neural Architectures, VAST 2019



TYPICAL ML PROCESS





DecisionTree? SVM?

kNN?

Regression?

K-means?







THE AUTOML APPROACH





AUTOML IS GREAT, BUT...

- The use of AutoML significantly reduces the effort of data modeling. However, AutoML still needs help to be useful!
- 1. Data Preparation
 - AutoML lacks "problem context" and "domain knowledge."
- 2. Query Specification
 - AutoML needs a carefully specified task. Specification is hard.
- 3. Model Comparison and Selection
 - AutoML will find models with the highest "score" but don't make sense (e.g. over or under fitted)





GOAL: BEST OF BOTH WORLDS

- Integration of AutoML with VIS:
 - Let AutoML handle all the machine learning
 - Frees the user to do what they do best



Hypothesis:
 AutoML+VIS >
 AutoML >
 Human



SNOWCAT: WORKFLOW



SNOWCAT

Snowcat

A visual analytics tool for exploratory model analysis

Tufts, Georgia Tech, Wisconsin-Madison



EVALUATION OF VIS + AUTOML

- Part of the DARPA D3M (Data-Driven Discovery of Models) program
 - Evaluation conducted by NIST
- Evaluation:
 - AutoML > human experts
 - Snowcat > AutoML

System	AutoML score	AutoML + VIS score
VA-1	20.2123081	21.20615752
VA-2	19.8853693	19.88951619
VA-3	19.8853693	20.58194965
VA-4	19.8853693	19.88951619
VA-5	19.8853693	19.88951619
VA-6	19.885369	273.1738095
Snowcat	19.8853693	9.440946857

Challenge 1



Challenge 2

System	AutoML score	AutoML + VIS score
VA-1	0.89707928	0.891516
VA-2	0.90055633	0.90055633
Snowcat	0.89707928	0.90472879

Results are classification accuracy, higher is better



Results are regression error (RMSE), lower is better



Gabriel Appleby (Tufts)



Anderson

(Novartis)

Erik



Alex Telea

(Utrecht)



Mateus Espadoto (Univ of Sao Paulo)

HUMAN-ML COLLABORATION: "UNPROJECTION"



PROJECTION (DIMENSION REDUCTION)

• Projection of n-dimensional space to 2D is the *bread and butter* of visual analytics

 There exist a number of techniques, both linear (PCA, MDS, Spectral) and non-linear (t-SNE, UMAP, Isomap, SOM).





PROJECTION AS A 1-WAY STREET

- Typically, projection is considered a 1-way operation.
 - That is, we can apply a projection function to a data point $P \in R^d$ to find its position in 2D (x, y)
 - However, given a 2D position (x, y), it is difficult to recover its high-dimensional position P
- But there might be a solution!
 - Since some projections learn a 2D manifold (e.g. UMAP), there is in fact a finite number of P's given an (x, y)





PROJECTION AS A 2-WAY STREET (1/3)

- Why this could be really useful:
 - Hypothesis Generation
 - Demo



- NSF project on material synthesis
 - Each dot is a known structure
 - Prediction of material structure can save time in unnecessary syntheses (expensive and time consuming)

vlalt

PROJECTION AS A 2-WAY STREET (2/3)

• Visualizing the manifold learned by a projection function (e.g. t-SNE)



(Middle): t-SNE projection of a 3D sphere is like "peeling an orange", but it's hard to see where the boundaries are

PROJECTION AS A 2-WAY STREET (3/3)

 Visualizing decision boundaries (of highdimensional data)



*(right): shows the agreement of an ensemble of 9 classification models. Dark blue and dark red backgrounds indicate agreement between all 9 models. Lighter hues represent the amount of disagreements

HOW IT WORKS: AUTOENCODER



Espadoto et al. "UnProjection: Leveraging Inverse-Projections for Visual Analytics of High-Dimensional Data." In submission to TVCG.

HUMAN+DATA SYSTEMS

- What do all good marriages have in common?
- Collaboration (Human+AI)
 - AutoML + VIS
 - Snowcat
 - Deep learning of projections
 - "UnProjection"

Communication (Human+DB)

- Computers and humans need to "talk" to each other quickly and effortlessly.
- Data systems need to respond to a user's interactions within 500ms
 - ForeCache, Kyrix
 - NeuralCubes









Wenbo Tao (MIT) Mike Stonebraker (MIT) Leilani Battle (Maryland)

DATA SYSTEM FOR SUPPORTING INTERACTIVE DATA EXPLORATION: DATABASE OPTIMIZATION

R. Chang et al., Dynamic Prefetching of Data Tiles for Interactive Visualization. SIGMOD 2016

R. Chang et al., Kyrix: Interactive Pan/Zoom Visualization at Scale. EuroVis 2019

R. Chang et al., Smile: A System to Support Machine Learning on EEG Data at Scale. VLDB 2019

R. Chang et al., Kyrix-S: Authoring Scalable Scatterplot Visualizations of Big Data. InfoVis 2020

INTERACTIVE EXPLORATION OF BIG DATA



Visualization on a Commodity Hardware



Large Data in a Data Warehouse



FORECACHE: PREDICTIVE PRE-FETCHING

- If we know what the user might do next, we can pre-compute the task and pre-fetch the result
 - Reduces user wait time (e.g. after they click a button to perform a task)







- ForeCache: Three-tiered architecture
 - Thin client (visualization)
 - Backend (array-based database)
 - Fat middleware



R. Chang et al., Dynamic Prefetching of Data Tiles for Interactive Visualization. SIGMOD 2016

PREDICTIONS, AN EXAMPLE

- What is the user going to do next?
 Likely "right" or "up"
- Prediction of behavior is pretty easy IF we know that the user is "navigating"
- But what if the user is doing something else?



PREDICTING USER ACTIONS

- Two-tiered approach using Markov
- First tier: predict what "phase" of analysis the user is in.
- Second tier: given a "phase", use phasespecific models to predict user's next actions.
 - For example, for "navigation", use momentum-based prediction
 - For others, use a combination of access-frequency, statistical distribution, SIFT (image-based), etc.



EVALUATION: PREDICTION ACCURACY

 Comparison against existing techniques

- "Random guessing" accuracy is: k/n
 - n: number of possible user actions
 - k: number of allowed "guesses"



KYRIX: GENERALIZED PAN-ZOOM

- Kyrix is an extension to ForeCache with two improvements:
- 1. It does not rely on pre-computed tiles
 - Data is stored in a modified R-tree
 - Update to an R-tree is much cheaper than recomputation of tiles
- 2. Kyrix includes a specification language
 - Authoring of a Kyrix "app" is fast and easy
 - The language, based on D3, allows for a wide range of visualization designs and layouts



$\langle SSV \rangle$::=	$\langle Marks \rangle \langle Layout \rangle \langle Data \rangle [Config]$	(1)
; marks			
(Marks)	::=	(Cluster)[Hover]	(2)
(Cluster)	::=	(Mode) (Aggregate) [Config]	(3)
(Hover)	::=	(Ranklist) (Boundary) [Config]	(4)
(Mode)	::=	Circle Contour heatmap	
		Radar Pie (Custom)	(5)
(Aggregate)	::=	$\langle Dimension \rangle * \langle Measure \rangle +$	(6)
(Ranklist)	::=	$\langle Topk \rangle (Tabular \mid \langle Custom \rangle)$	(7)
(Boundary)	::=	Convex Hull BBox	(8)
(Custom)	::=	Custom JS mark renderer	(9)
(Dimension)	::=	(Field)[Domain]	(10)
(Measure)	::=	<pre>{Field> (Function> [Extent]</pre>	(11)
(Topk)	::=	A positive integer	(12)
(Domain)	::=	A list of string values	(13)
(Function)	::=	Count Sum Avg Min	
		Max Sqrsum	(14)
; layout			
(Layout)	::=	$\langle X \rangle \langle Y \rangle \langle Z \rangle$ [Theta]	(15)
$\langle X \rangle$::=	(Field)[Extent]	(16)
$\langle \mathbf{Y} \rangle$::=	(Field)[Extent]	(17)
$\langle Z \rangle$::=	(Field) (Order)	(18)
(Theta)	::=	A number between 0 and 1	(19)
(Field)	::=	A database column name	(20)
(Extent)	::=	A pair of float numbers	(21)
(Order)	::=	Ascending Descending	(22)
; data			
$\langle Data \rangle$::=	a database query	(23)
; config			
(Config)	::=	Key value pairs	(24)



KYRIX: EXAMPLE APPS (NBA)



Data: All plays in the 2017-2018 NBA season with 560k records

KYRIX: SPECIFICATION LANGUAGE

```
var logoTransform = new Transform("logoTransform",
1
       "select * from teams;",
2
3
       "nba",
      function (row){
4
5
          var id = parseInt(row[0]);
6
          var y = Math.floor(id / 6);
          var x = id - y * 6;
8
          var ret = [];
          9
          ret.push((y * 2 + 1) * 80 + 100); // y coordinate of logos
10
          for (var i = 1; i \le 4; i ++)
11
                               // raw data attributes
12
              ret.push(row[i]);
13
          return ret;
14
      },
      ["x", "y", "team_id", "city", "name", "abbr"]);
15
```

1 var selector = function (row) { 2 return true; 3 3; var viewport = function (row) { 4 5 return [0, 0]; 6 }; var predicate = function (row) { 7 8 return { 9 "layer 0" : "home_team=" + row.team_id + " and " 10 + "away_team=" + row.team_id, 11 "layer 1" : "id=" + row.team_id 12 3; 13 }; p.addZoom(new Zoom(logoCanvas, 14 15 timelineCanvas, selector, 16 17 viewport.

// ====== teamlogo -> teamtimeline =======

Transition Functions

18

predicate);

7

KYRIX: SYSTEM ARCHITECTURE (REALTIME)





KYRIX: EXAMPLE APP (MGH EEG)



Dataset: 17 million patients (shown as dots in the scatterplot) along with the 100+ million data points for all of their EEG data (right).



EVALUATION

- Benchmark Performance
 - Beats target of 500ms by 10x



- User study
 - 8 participants (4 with visualization programming experience, 4 without)
 - Participants asked to recreate (parts of) the NBA example*
 - Completion time of avg: 26.25 minutes, stdev: 7.07 minutes

*Participants wrote the placement functions and the transition functions of the top two layers, but did not have to write the transform or the rendering functions.





Zhe Wang (Arizona)



Dylan Cashman (Tufts)



Mingwei Li (Arizona)



Jixian Li (Arizona)



Matt Berger (Vanderbilt) Josh Levine (Arizona)



Carlos Scheidegger (Arizona)

[v]alt

DATA SYSTEM FOR SUPPORTING INTERACTIVE DATA EXPLORATION: **NEURALCUBES**

INTERACTIVE EXPLORATION OF BIG DATA



Visualization on a Commodity Hardware



Large Data in a Data Warehouse



visit and the second se

NEURALCUBES -- BIG PICTURE IDEA



Source: US Bureau of Labor Statistics, Current Employment Statistics and TIP Strategie



Source: US Bureau of Labor Statistics, Current Employment Statistics and TIP Strategies

[v]alt

SAMPLE RESULTS



STORAGE SAVING

DataSet	Raw Data Size	# queries	# net paramy	network size	# iteration	training time	testing error
	100,000 (3.9MB)	1000k	50,653	232KB	1000	3 h 5 mins	0.022
	100.000 (3.9MB)	2000k	53,903	245KB	1000	3 h 8 mins	0.028
SPLOM	100,000 (3.9MB)	3000k	57,153	258KB	1000	8 h 44 mins	0.038
	100,000 (3.9MB)	4000k	60,403	271KB	1000	10 h 32 mins	0.048
	100,000 (3.9MB)	5000k	63,653	284KB	1000	10 h 17 mins	0.057
BrightKite NYC	79,357 (3.1MB)	8,860k	169,032	688KB	500	4 h 29 mins	0.124
BrightKite Austin	22,171 (868KB)	8,860k	169,032	688KB	300	5 h 55 mins	0.094
Flights (count)	5,092,321 (204MD)	856,689	24,267	152KB	250	1 h 16 mins	0.042
Flights (delay)	5,013,088 (175MB)	670,000	123,330	830KB	500	2 h 14 mins	0.025

- When original data is small, not much saving
- When original data is larger, the saving can start to become more significant

BUT WAIT, THERE'S MORE!!

Raw Data Size	# States	Model Size	Testing RAE
12k (1.8MB)	30k	798KB	3.70%
120k (18MB)	30k	798KB	2.04%
1.2m (180MB)	30k	798KB	1.35%
12m (1.8GB)	30k	798KB	0.97% 🔶

- As the raw data size goes up, so does the amount of training data
 - As a result, the amount of error actually decreases (while the model size stays constant)!!
 - In this case, the space saving is 2000X
 - (2 GB -> 1MB)

2D AUTOENCODER



RECAP: DATA SYSTEMS

- Human + ML:
 - SnowCat
 - AutoML + VIS system that is better than AutoML and human expert
 - "Unproject"
 - Learns the inverse of a projection function
- Human + DB:
 - ForeCach & Kyrix
 - Optimized database systems for pan-zoom visualization systems
 - NeuralCubes
 - Use of deep learning to compress the data







THE FUTURE OF DATA SYSTEMS FOR VIS



- Future of data analysis (and data science) will be at the intersection of VIS/HCI, Stats/ML/AI, and DB/Cloud
- Each community is becoming more aware of this need
 - KDD: IDEA workshop
 - SIGMOD: HILDA workshop
 - VIS: DSIA workshop
- Work I presented today is just scratching the surface.
 - How to integrate all three is hard
 - Space is wide open. Join us!





QUESTIONS?

REMCO @CS.TUFTS.EDU



