Integrating Predictive Models with Interactive Visualization



Jian Zhao, Ph.D., Assistant Professor Cheriton School of Computer Science University of Waterloo <u>www.jeffjianzhao.com</u> jianzhao@uwaterloo.ca

Short bio		
	Researcher @ Autodesk, Toronto	Assistant Professor @ U Waterloo
	2015	2019
2009	2016	
Ph.D. @ U Toronto	Researcher @ FXPAL, Palo Alte	0

Al Machines

a 1010111110100 m

1010010011

0010101

Data

110110

01001101100

1010.01

*0111011011010101011 *01000001100100 101010110011110* 001010001100* 001010000* 001010000* 00100000* 00100000* 0010000* 0010000* 0010000* 0010000* 0010000* 0010000* 0010000* 0010000* 001000* 0010000* 0010000* 0010000* 0010000* 0010000* 00* 000* 00

01011

Humans

All continuously growing fast!

I investigate advanced visualizations (vis) that promote the interplay among data, machines (models), and humans (users) in real-world data science applications.



"My input data looks similar, but my classifier performs quite different... Why?"

Х	Mean	:	54.26
Y	Mean	:	47.83
Х	SD	:	16.76
Y	SD	:	26.93
Co	orr.	:	-0.06

Bella, Data Scientist



Matejka et al, Same Stats, Different Graphs: Generating Datasets with Varied Appearance and Identical Statistics through Simulated Annealing, CHI'17

"I'm building a neural network classifier. I tried many ways, but it doesn't work... Why?"



Bella, Data Scientist





Tensor Flow Playground, http://playground.tensorflow.org/

"I finally got some good results, but my boss couldn't understand them..."



Bella, Data Scientist

Algorithm	Mean Absolute Error	Root Mean Squared Error	Relative Absolute Error	Relative Squared Error	Coefficient of Determination
Bayesian Linear Regression	3.258579	4.630709	0.508712	0.279148	0.720852
Neural Network Regression	2.738169	3.64928	0.427469	0.173362	0.826638
Boosted Decision Tree Regression	2.091976	2.84549	0.326588	0.105403	0.894597
Linear Regression	3.423023	4.651562	0.534385	0.281668	0.718332
Decision Forest Regression	2.439228	3.497692	0.3808	0.159258	0.840742





Visualization is critical in data analysis workflow



Top machine learning and data science methods used at work



http://businessoverbroadway.com/top-machine-learning-and-data-science-methods-used-at-work

Creating effective visualizations is hard

Problem/domain specific No easy one-size-fits-all solution

Technical skills Matplotlib, D3.js, ggplot2, ...

Sense of design Huge design space











Make sense of results

Data analysts General users



Make sense of data

Make sense of models

Make sense of results

Explore complex data with visualization recommendations

Comprehend missing link prediction in bipartite networks Leverage video recommendations in online learning







ChartSeer

MissBiN





Make sense of data



Exploring large information space



Challenges

Continuously making decision in a large parameter space Which data variables to explore? What kind of charts to use?

Lacking a holistic view of the analysis space How is the current status? Where am I?

Exploring large information space with recommendation





J. Zhao, M. Fan, M. Feng, ChartSeer: Interactive Steering Exploratory Visual Analysis with Machine Intelligence, TVCG



System architecture



Chart summarization



#18-11

Average Faculty Salary

Controlled user study

Between-subjects design 24 participants (13 females and 11 males)

Interface conditions ChartSeer v.s. Baseline

Dataset US college statistics (18 variables)

Tasks

Summarization task Exploration task

Inspection Pa	net	Data and Charts Panel				
<u>Г</u>	af98 a	Name	* Miles,per,Gallon ()	Cylinders (Displacement	Horsep
1	000000	amc ambassador brougham	13		360	175
1.0	2.000	amc ambassador dpl	15		390	190
1		amc ambassador sat	17		304	150
1 1	· · · · · · · · · · · · · · · · · · ·	amc concord	19.4	6	232	90
- 4	00 0000 00	amc concord	24.3	4	151	90
		amc concord d/l	18.1	6	258	120
0	5 5 5 5 5	amo concord di	23	4	151	null
	Acceleration	amc concord dl 6	20.2	6	232	90
Add	Update Remove attent	amc gramin	21	6	199	90
Mark	aciet	amo gremin	19	6	232	100
x	Acceleration a -	amo gremin	18	6		100
¥	Horsepower	amc gramin	20	6	232	100
Collor	· • •	amp homet	18	6	199	97
Size	- 0 0	amp homet	18	6		100
Shape	· 🔁	amp homet	19	6	232	100
Preview	Cancel	amp homet	22.5	6	232	90
		amc homet sportabout (sw)	18	6	258	110
1 - (amo matador	18	6	232	100
Z -	"encoding": {	amo matador	14		304	150
4	"field": "Horsepower",	amo matador	16	6	258	110
5	"type": "quantitative"	amo matador	15	6	258	110
6	h	amo matador	15.5		304	120
8	"field": "Acceleration",	amc matador (sw)	15		304	150
9	"type": "quantitative"	amc matador (sw)	14		304	150
10	1	amo pacer	19	6	232	90
12	"mork": "point"	amc paper d/l	17.5	6	258	95
13 }		amp rebel ant	16		304	150
		amc rebel ant (rw)	null		360	175
		amc spirit di	27.4	4		80
		audi 100 la	24	4		90
		audi 100is	20	4	114	91
		audi 100is	23	4	115	95
		audi 4000	34.3	4	97	78
		aud: 5000	20.3	5	131	103
		aud 5000s (dese)	36.4	5		67
Preview	Cancel					

Results of user behaviors

Participants added more charts but updated less charts using ChartSeer

ChartSeer led to a broader range of data variables and visual encodings

ChartSeer encouraged more focused exploration of data variables

ChartSeer allowed for data exploration from more heterogenous visual perspectives



🔄 ChartSeer



Questionnaire results









Make sense of models

	Name	e State	Gender	Score	
	0 Georg	ge Arizona	м		63
1	1 Andre	ea Georgia	F		48
	2 mich	eal Newyork	M		56
	3 magg	ie Indiana	F		75
	4 Ravi	Florida	M	NaN	
	5 Xien	California	M		77
	6 Jalpa	NaN	NaN	NaN	
	7 NaN	NaN	NaN	NaN	
		X	\frown		\sim
	0)		
	(6		1
0)				(
/ `				5	1
1	7			X	-
(0	1		
5	-		6		
)		20)	(0
				1	

×.

 \sim

"Missing" links in bipartite networks







Missing link prediction



C - 5: 0.974 D - 2: 0.965 E - 1: 0.873 B - 3: 0.852

• • •

Analysts' questions



MissBiN



J. Zhao, M. Sun, F. Chen, P, Chiu, MissBiN: Visual Analysis of Missing Links in Bipartite Networks, VIS'19 J. Zhao, M. Sun, F. Chen, P, Chiu, Understanding Missing Links in Bipartite Networks with MissBiN, TVCG

Addressing the questions with MissBiN



A missing link prediction algorithm

An interactive visualization

A comparative analysis approach

Prediction of missing links

1. Predict the missing links with standard methods (e.g., common neighbors [Chang12])

2. Discover all maximal bicliques, complete subgraphs, of the network (e.g., using MBEA [Zhang14])

3. Re-rank the missing links based on the overlap of bicliques

Algorithm 1: Missing link ranking. Input : A list of bi-cliques, $L = \{C_i = (X_i, Y_i, E_i)\}$, detected in a bipartite network G = (X, Y, E). Output: Weights, w, for all non-observed (missing) links in G. foreach $e \in X \times Y - E$ do $w_e \leftarrow 0;$ 3 end 4 foreach bi-clique pair (C_i, C_j) from L do $o \leftarrow \frac{|X_i \cap X_j|}{|X_i \cap X_j|}$ if o < threshold the continue; end $\frac{|X_i \cap X_j| \cdot |Y_i \cap Y_j|}{|X_i - X_j| \cdot |Y_j - Y_i| + |X_j - X_i| \cdot |Y_i - Y_j|};$ 9 foreach $e \in \{\langle x, y \rangle; x \in X_i - X_j, y \in$ $Y_j - Y_i \} \cup \{(x, y); x \in X_j - X_i, y \in Y_j - Y_i\}$ do $w_e \leftarrow w_e + w_i$ 11 end end

In step3, for each pair of bicliques, ...



Re-ranking predicted missing links

Weights computed in step3, based on bicliques information

$$s'_e = w_e \cdot s_e$$

Scores computed in step1, based on standard methods

Evaluation of missing link prediction

Test on 3 datasets

Person-place network from Atlantic Storm corpus [Hughes05] User-conversation network from Slack group communication

Compare with 5 base methods Jaccard coefficient (JA) common neighbors (CN) Adamic-Adar coefficient (AA) preferential attachment (PA) random walk (RW)

Link prediction results

Mostly, PA has the largest performance gain Secondly, CN performs well Original method Our method ^performance gain

		/	Atlantic S	Storm	Ι	Slack	comm	unication	Ι	VI	S public	ations	1	Α	tlantic S	storm*	Ι	Slack	commu	nication*
R-Precision	JA CN AA PA RW	.395 .359 .440 .021 .467	.428 .580 .455 .585 .531	.032 .221 .016 .564 .064		.421 .179 .202 .025 .451	.424 .424 .219 .590 .464	.004 .245 .016 .565 .013		.105 .066 .107 .000 .096	.126 .129 .152 .012 .112	.020 .063 .045 .012 .016		.405 .342 .405 .007 .487	.423 .552 .418 .226 .552	.018 .209 .013 .219 .065		.332 .202 .203 .096 .398	.341 .339 .218 .432 .407	.009 .138 .015 .336 .009
AUC PR	JA CN AA PA RW	.398 .305 .398 .008 .435	.451 .607 .429 .566 .516	.053 .302 .031 .557 .082		.301 .148 .170 .021 .346	.325 .326 .200 .528 .377	.024 .178 .030 .507 .031		.039 .025 .039 .000 .039	.053 .063 .068 .003 .052	.014 .039 .029 .002 .013		.393 .300 .379 .005 .448	.444 .578 .406 .212 .525	.051 .277 .027 .207 .078		.225 .133 .153 .055 .279	.249 .285 .177 .406 .308	.024 .152 .025 .352 .029

Jaccard coefficient (JA), common neighbors (CN), Adamic-Adar coefficient (AA), preferential attachment (PA), random walk (RW)

Addressing the questions with MissBiN



A missing link prediction algorithm

An interactive visualization

A comparative analysis approach



Showing 1 to 153 of 153 entries

changed metric

Link detection	Numic. +	Score filter	Matrix order	by id	*	Compute Motifs and Metrics



node 10	degree (before)	degree (after)	closenese (before)	closenses (after)	betweenness (before)	betweenness (after)
US7	0.016088	0.0000	6.33609	6.0000	0.0000	0.0000
1.58	0.19672	0.0000	0.41605	6.0000	0.011734	0.0000
UBD	0.11475	0.0000	0.44727	£.0000	0.000616	0.0000
087	0.12500	0.0000	0.01007	6.0080	0.027018	0.00000
C136	0.12500	0.0000	0.56296	6.0000	0.0006435	6.00000
0140	0.12500	0.0000	0.56296	6.0000	0.0006435	6.0000
C166	0.19625	0.0000	0.63866	0.0000	0.0400N3	0.0000
Showing 1 to 90 of	10 estries					

Search:

Link detection	bi_adamic, +	Score filter	Matrix order	by bital-rum +	Compute Moths and Metric



8.43601

6.0008

0.0048347

0.0000

0764

Showing 1 to 90 of 90 entries

0.12500

0.0000

Evaluation of MissBiN

Interview study

A management school professor on exploring organizational communication networks

A computer scientist on investigating relationships of crimes and locations in Washington DC

Case study The Sign of the Crescent [Hughes03] 41 fictional intelligence reports Extracted person-location network 49 persons and 104 locations, with 328 links Analysis task

Identify suspicious persons and activities from the reports

Make sense of results

Exploring large information space with recommendation





Linear ranked list is not enough

Semantic map significantly improves users' comprehension capability compared to a ranked list [Peltonen 2017]

Orienteering helps understand and trust the answers using both prior and contextual information [Teevan 2004]

Support stepping behavior by clustering the information or suggesting query refinements [Teevan 2004]

Mike, the confused



Want to solve an optimization problem in his work Just watched #19 - choosing stepsize and convergence criteria



Recommendations:

- 1. Sparse models selection
- 2. Dirichlet distribution
- 3. Gradient descent intuition
- 4. Hill climbing

5.



J. Zhao, C. Bhatt, M. Cooper, D. Shamma, Flexible Learning with Semantic Visual Exploration and Sequence-Based Recommendation of MOOC Videos, CHI'18





Convergence criteria

For convex functions, optimum occurs when

àg(-) - 0

In practice, stop when





Algorithm:





System architecture



Recommendation engine

Content-based recommendation Based on TF-IDF

Sequence-based re-ranking Topic similarity score (TS) Global sequence score (GS) Local sequence score (LS)

Sub-sequence aggregation Greedy search down the ranked list

Dataset ~4000 videos, ~350 hours running time, from Coursera, EdX, and Udacity

$$egin{aligned} \hat{S}_{TS}(P^{(q)},P^{(r)}) = \mathrm{J}(V_q,V_r) + rac{1}{Z}\sum_{z=1}^Z \delta\left(|P^{(q)}(z) - P^{(r)}(z)| \leq t
ight) \ &S_{GS}(V_q,V_r) = rac{1}{N}\sum_{i=0}^N c_i rac{y_i}{|V_r|} + rac{s_i}{D_G} \ &S_{LS}(V_q,V_r) = rac{1}{M}\sum_{i=0}^M c_i rac{y_i}{|V_r|} + rac{s_i}{D_q} \end{aligned}$$

Visualization generation

Multidimensional scaling (MDS) in feature space Rotate to comply with left-right browsing flow Tune positions to avoid overlap Merge consecutive videos

Hierarchical clustering Context-based region division Voronoi tessellation

Topical keywords extraction Force-directed placement



Scenario I: "I missed anything?"



Mike

Confused about this lecture. Wants to check if missed anything.



Scenario II: "I want to know more."



Lisa

Already knows about this. Wants to extend her horizon.



Used by MOOC instructors

Semi-structured interviews with two university instructors

"I normally don't look at what others teach, but the tool provides the awareness of related lectures, so I could borrow some materials to enhance my lecture, and avoid unnecessary duplication." "If you see one lecture is here [on the Exploration Canvas], then you go very far for the second lecture, and back here again for the third lecture, you should really think about reordering the content presented in the videos." One more thing...

Thank all my collaborators!





Available on https://www.jeffjianzhao.com/webapp/EgoLines/egolines.html

Another thing...

Welcome to apply to Waterloo HCI

WATERLOOHC CHERITON SCHOOL OF COMPUTER SCIENCE

http://hci.cs.uwaterloo.ca/



Jeff Avery Lecturer



Géry Casiez Adjunct Professor



Keiko

Katsuragawa

Research Assistant Professor and Research Officer at the NRC

ISIEZ fessor



Edward Lank



Edith Law Assistant Professor







Jian Zhao Assistant Professor



Integrating Predictive Models with Interactive Visualization

Jian Zhao, *Ph.D.*, Assistant Professor Cheriton School of Computer Science University of Waterloo <u>www.jeffjianzhao.com</u> <u>jianzhao@uwaterloo.ca</u>