NYU 2015 Conference on Digital Big Data, Smart Life, Mobile Marketing Analytics October 23, 2015

Using Topological Data Analysis to Explore **Emergent Consumer Experience** from Digital Interactions Tom Novak

Donna Hoffman

The George Washington University Center for the Connected Consumer

Interactivity is Evolving and New Consumer Experiences are Emerging

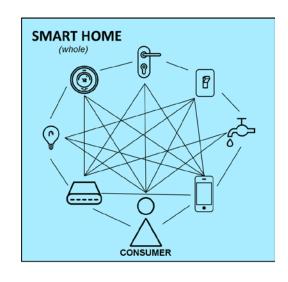
The consumer IoT represents a multitude of on-going, evolving heterogeneous interactions among many different components: C2M, M2M, M2P, C2C – a set of recurrent "assemblages" (Hoffman and Novak 2015).

What kind of insights can we derive about emergent consumer experience in the IoT (the "possibility space") from actual interactions?

These interactions represent a lot of digital "big data" – very high dimensionality data consisting of often millions of ongoing interactions among complex, heterogeneous component devices (and consumers!).

So IoT big data are big in two ways: 1) size, and more importantly, 2) complexity.

The challenge: Get from the complex unknown to the more interpretable known.



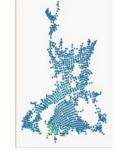
Making Sense of of Digital Big Data from the IoT

Predictive Analytics approach: Fit predictive models to the data. But the complexity of the data means hypothesis testing is often challenging. We need to know what questions to ask. Are we asking the right questions? With big data, insights can be slow.

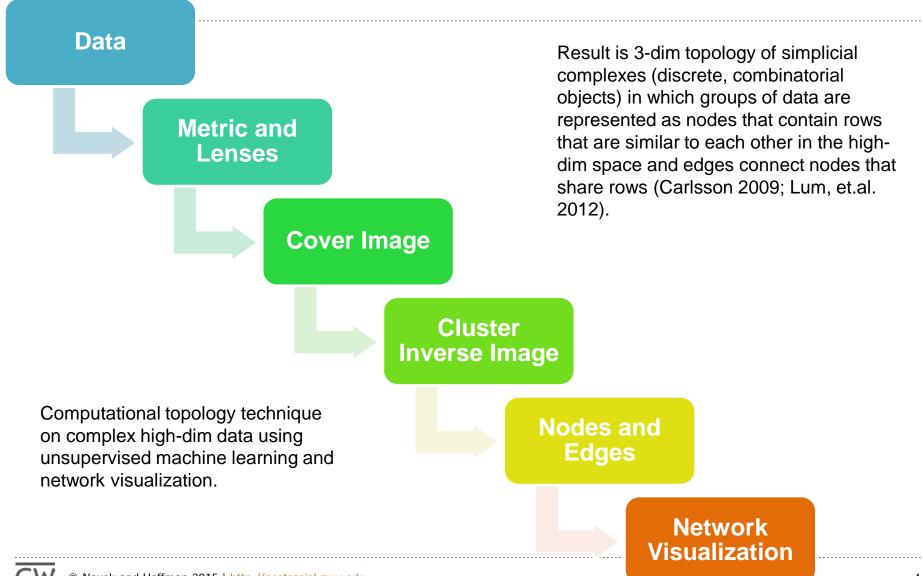
Conventional approaches for reduction and visualization: Use linear and nonlinear dimension reduction techniques such as PCA, MCA, and MDS. But, even if they work, are sensitive to distance metrics and do not preserve topological structures of the data.

Data-Driven Discovery Approach: Hypothesis-free approach based on computational topology to qualitatively analyze functions on very high-dimensional data and visualize the data structure in low-dimensional topological spaces. Topological data analysis (TDA) reveals structures in the data that have invariant properties and can propel insight and improve hypothesis-generation and predictive modeling; "digital serendipity" (Singh 2013).

132 million data points _



TDA Framework for Generating Topological Networks



Data: 120,253 IFTTT (If This Then That) Rules

"Turn on my lights when I get close to my home"

(Title Words)

563 binary variables

IF

IOS Location

(Trigger Channel)

86 binary variables

You enter an area

(Trigger Words)

280 binary variables

THEN

Phillips Hue

(Action Channel)

69 binary variables

Turn on the lights

(Action Words)

103 binary variables

Data Description

120,253 IFTTT rules were created by 60,230 IFTTT users over a 3 year period from mid 2011 to mid 2014. (8404 unique rules)

68% of users created 1 rule
22% of users created 2-3 rules
8% of users created 4-9 rules
2% of users created 10+ rules
(23,496 total rules or 20% of data)

1101 binary variables for trigger channels, action channels, trigger words, action words & title words.

132 million data points



Most Popular IFTTT Trigger and Action Channels

IFTTT Trigger Channel	IFTTT Action Channel		
Feed (RSS feeds) 25%	Twitter 15%		
Instagram 9%	Email 11%		
Date & Time 8%	SMS 8%		
Weather 8%	Evernote 7%		
Facebook 5%	Facebook 7%		
Gmail 4%	Dropbox 6%		
YouTube 2%	Facebook Pages 4%		
Twitter 2%	Google Drive 4%		
WordPress 2%	Tumblr 3%		
Tumblr 2%	Pocket 3%		
Total top trigger channels 68%	Top top ten action channels 74%		

Some IFTTT Rules Were Reinvented Again and Again

Trigger Channel	Trigger	Action Channel	Action	# of users creating this rule
Feed	New feed item	Email	Send me an email	4987
Feed	New feed item	Twitter	Post a tweet	2868
Date & Time	Every hour at	Twitter	Post a tweet	1114
Facebook	You are tagged in a photo	Dropbox	Add file from URL	926
Weather	Sunrise	SMS	Send me an SMS	416

8404 unique rules, where:

39% of unique rules were created by 1 user 24% of unique rules were created by 2-3 users 18% of unique rules were created by 4-9 users 19% of unique rules were created by 10+ users

1% of unique rules were created by 200+ users (49,201 total rules of 41% of data)

///////

Preliminary Research Questions

- 1. What IFTTT rules are people creating?
- 2. Are there interesting and important patterns that underlie the IFTTT rules that have been created?
- 3. Do these patterns suggest emergent themes?

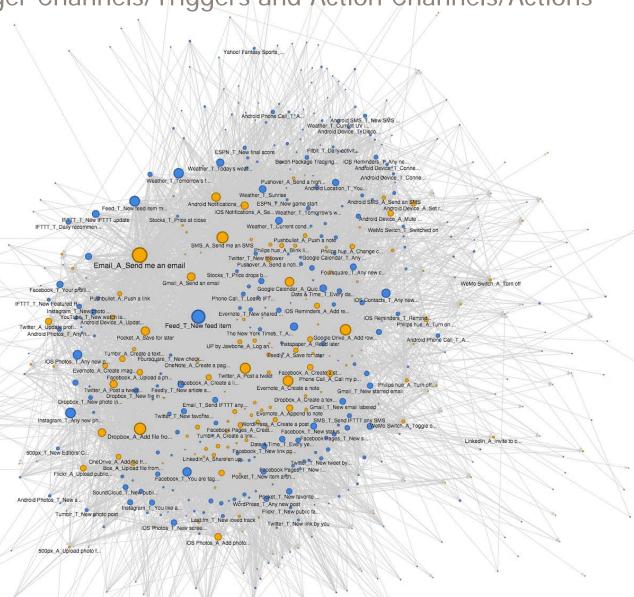
→ Let's take a look at some conventional reduction and visualization approaches first...



Network Graph of 553 Nodes for 120,253 IFTTT Rules Using All IFTTT Trigger Channels/Triggers and Action Channels/Actions

Network graph cannot reveal clear patterns in these complex data.

Only option would be to significantly reduce the number of nodes shown.



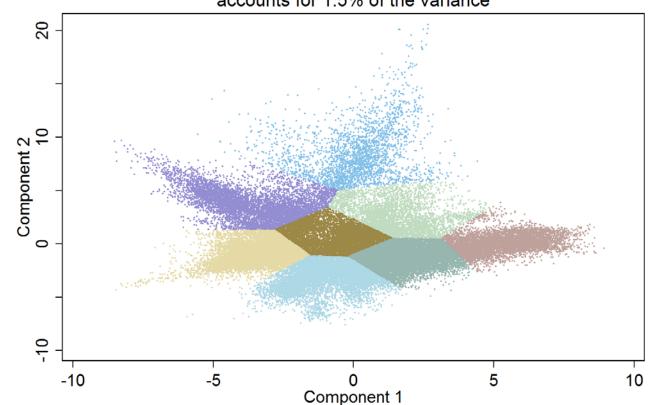
PCA and K-Means Clustering on the Component Scores

Some general features, but complexity is not revealed:

Component 1 identifies groups of rules about photos (+) versus email and SMS (-)

Component 2 separates weather and smart home rules (+) and new rules and feeds (-)

Plot of IFTTT Recipes on the First Two Components accounts for 1.5% of the variance



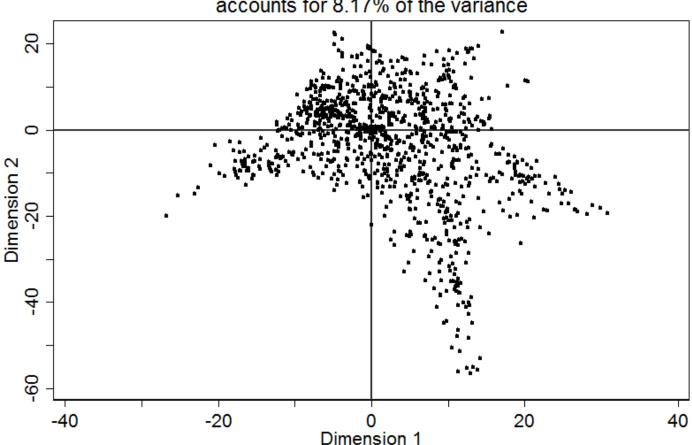


Multiple Correspondence Analysis of the Burt Matrix

MCA of IFTTT Recipes on the First Two Dimensions accounts for 8.17% of the variance

Some gross features, but no complexity:

dimension 1 separates weather (+) and photos (1) while dimension 2 contrasts reading (+) and the smart home devices and Android devices (-).





TDA Provides a Different Approach

Conventional dimension reduction and visualization approaches have difficulty revealing clear patterns in these complex data of 1100+ vars.

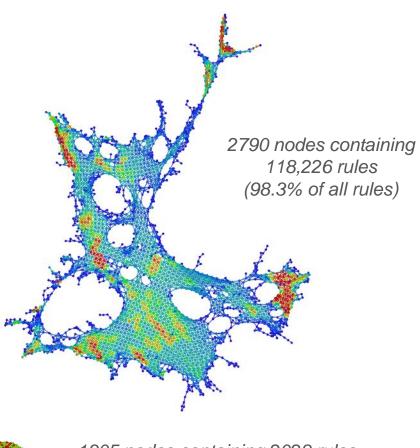
Need about 400 dimensions to account for most of the variance in the data. Took Stata 6 hours to compute the MCA solution.

General features that grossly separate rules are apparent, but it's difficult to see the more subtle behavioral patterns underlying rule creation, let alone potential emergent themes.

Topological data analysis offers an approach to visualization of network structure that is more revealing and useful.



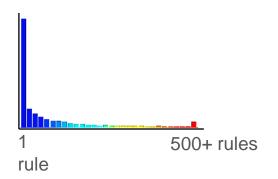
Ayadsi* TDA Solution of All 120,253 IFTTT Rules



1305 nodes containing 2028 rules (1.7% of all rules) Each node is a cluster of rules.

Nodes connect if they have rules in common.

Color indicates the number of rules in each node:



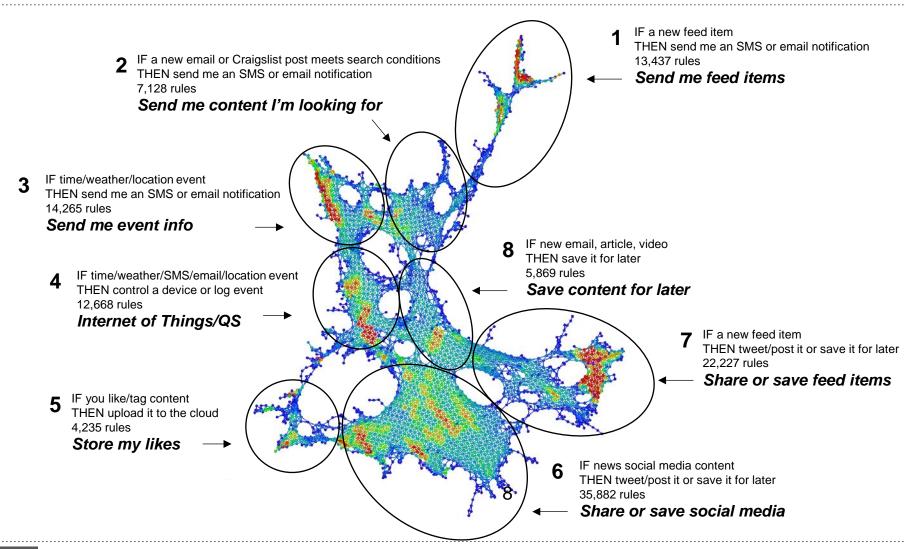
Metric: Hamming

Lenses: MDS 1 and MDS 2 (resolution 60, gain 1.6)

*We used the Ayasdi 3.0 software platform (ayasdi.com, Ayasdi, Inc., Menlo Park, CA) to perform the TDA on the IFTTT data



TDA Reveals the Possibility Space of IFTTT Rules



Using Data Lenses to Include Known Structure

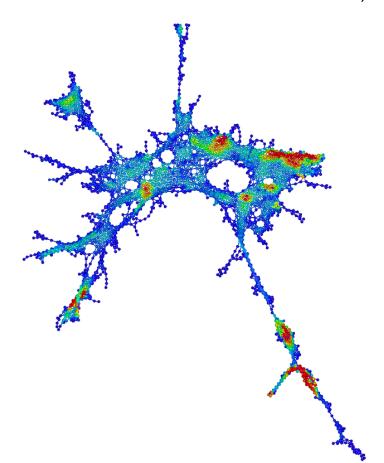
We can use categorical variables as "data lenses" to separate the IFTTT rules into groups with distinct structures. For example:

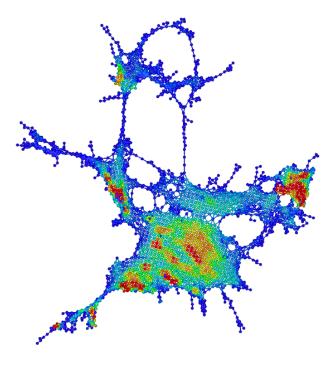
- Whether the IFTTT rule did or did not include a social media channel (e.g. Facebook, Twitter, YouTube, etc.). Does the structure of IFTTT rules vary based on whether the rule includes a social media channel?
- The year the IFTTT rule was created. Does the structure of IFTTT rules change over time?

Social and Non-Social IFTTT Rules

54,300 Non-Social Rules (45% of rules do not use a Social Media channel)

65,953 Social Rules (55% of rule use a Social Media channel)



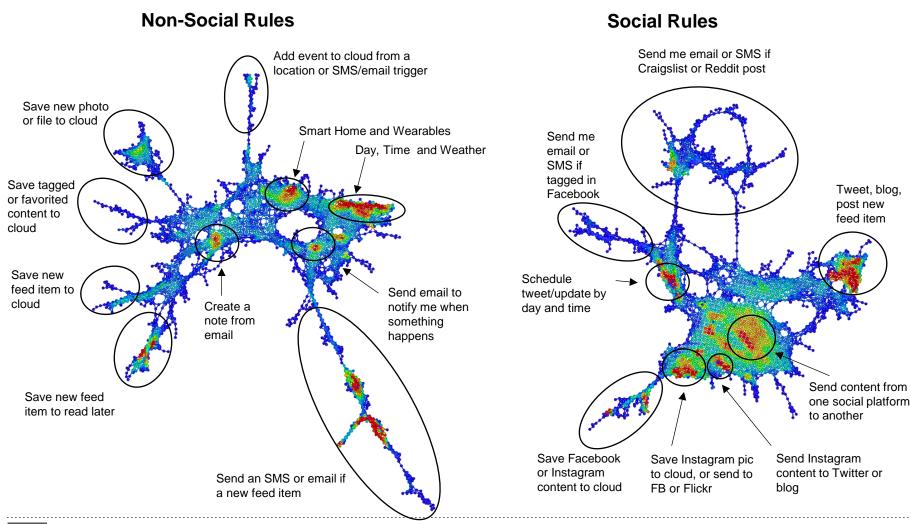


Metric: Hamming

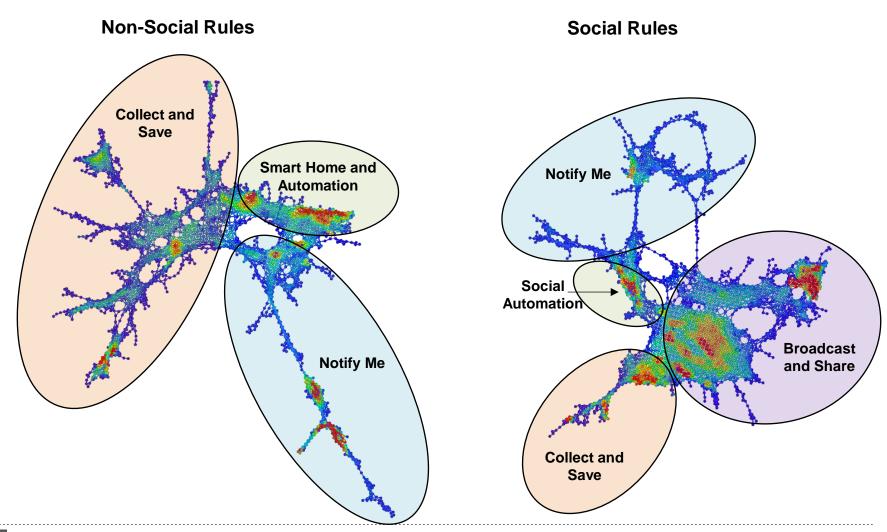
Lens: MDS 1 & 2 (resolution 60, gain 1.9)

Data Lens: Social Media (2 groups)

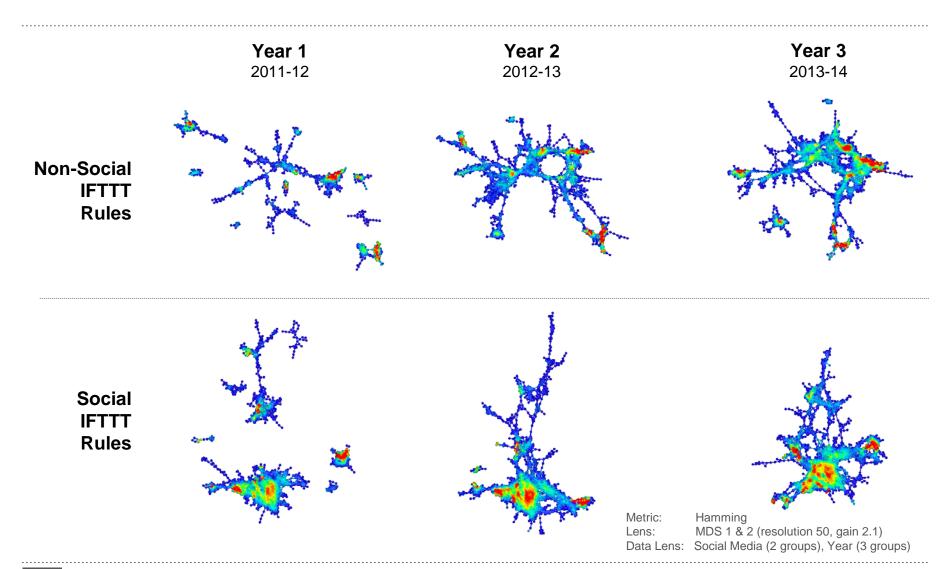
Social and Non-Social IFTTT Rules



Social and Non-Social IFTTT Rules



TDA of IFTTT Rules Over Time



Collect & Save

Automation

Notify Me

Broadcast & Share

The basic structure of IFTTT emerged in year 1, but became more organized and interconnected in years 2 and 3.

Year 1 Year 2 Year 3 2011-12 2013-14 2012-13 **Non-Social IFTTT Rules** Social **IFTTT Rules**

Summary

TDA provides a way to visualize what emerges from the IoT using interaction events as the unit of analysis.

IFTTT is an assemblage that emerges from the capacities of the components (i.e. triggers and actions) exercised in the interaction over time between apps and devices that are connected through the individual rules.

Based on our assemblage theory framework (Hoffman and Novak 2015), the topology represents the possibility space (DeLanda 2006, 2011) underlying the potential capacities of the IFTTT assemblage.

Supports productive hypothesis generation and subsequent predictive modeling.



The Ayasdi 3.0 software platform for topological data analysis (ayasdi.com) was used to construct all networks of the IFTTT data. The authors acknowledge the support of Devi Ramanan, Global Head – Product Collaborations, Ayasdi Inc., Menlo Park, CA

IFTTT data were collected from a crawl during May - June 2014 and are used with permission of IFTTT.com, San Francisco, CA.