



# Data Aggregation/Summarization and Business Implications

Dai Yao

Associate Professor of Marketing

Department of Management and Marketing



[dai@yaod.ai](mailto:dai@yaod.ai)

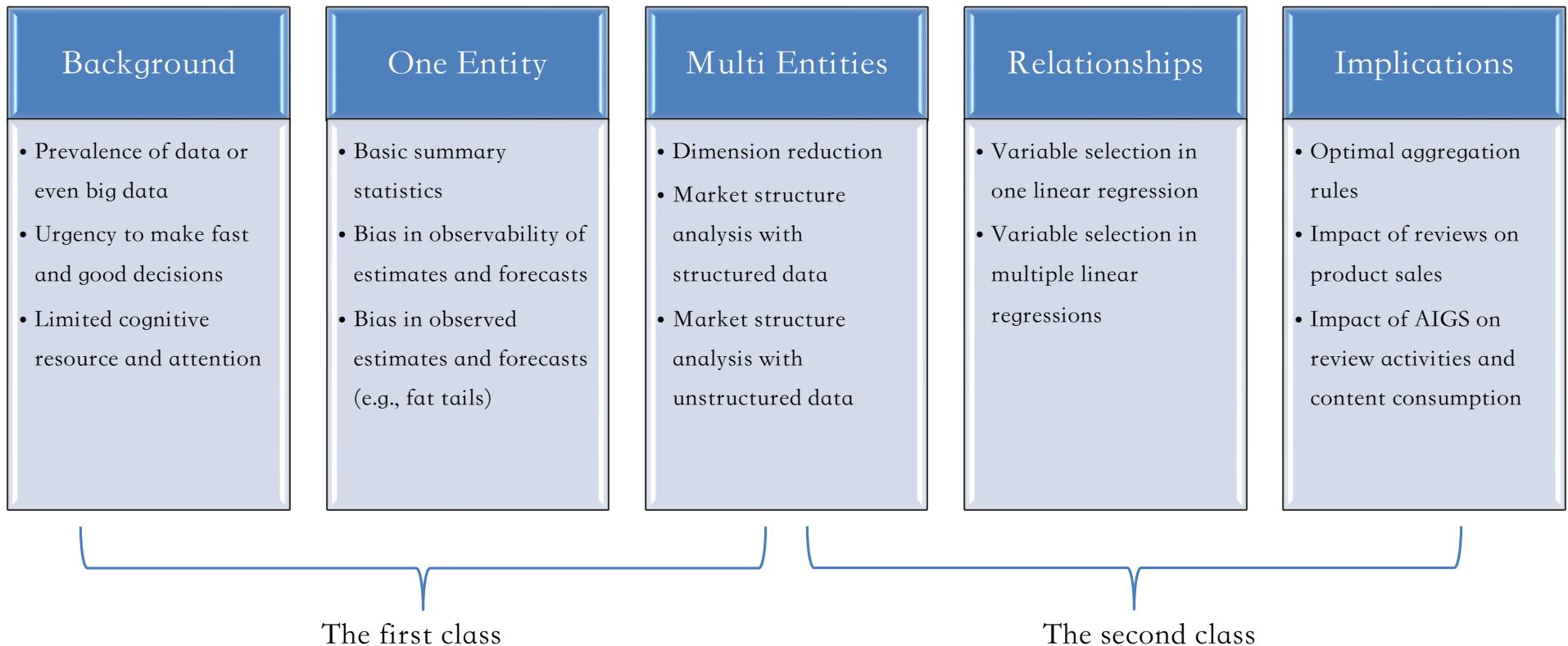


<http://www.yaod.ai>

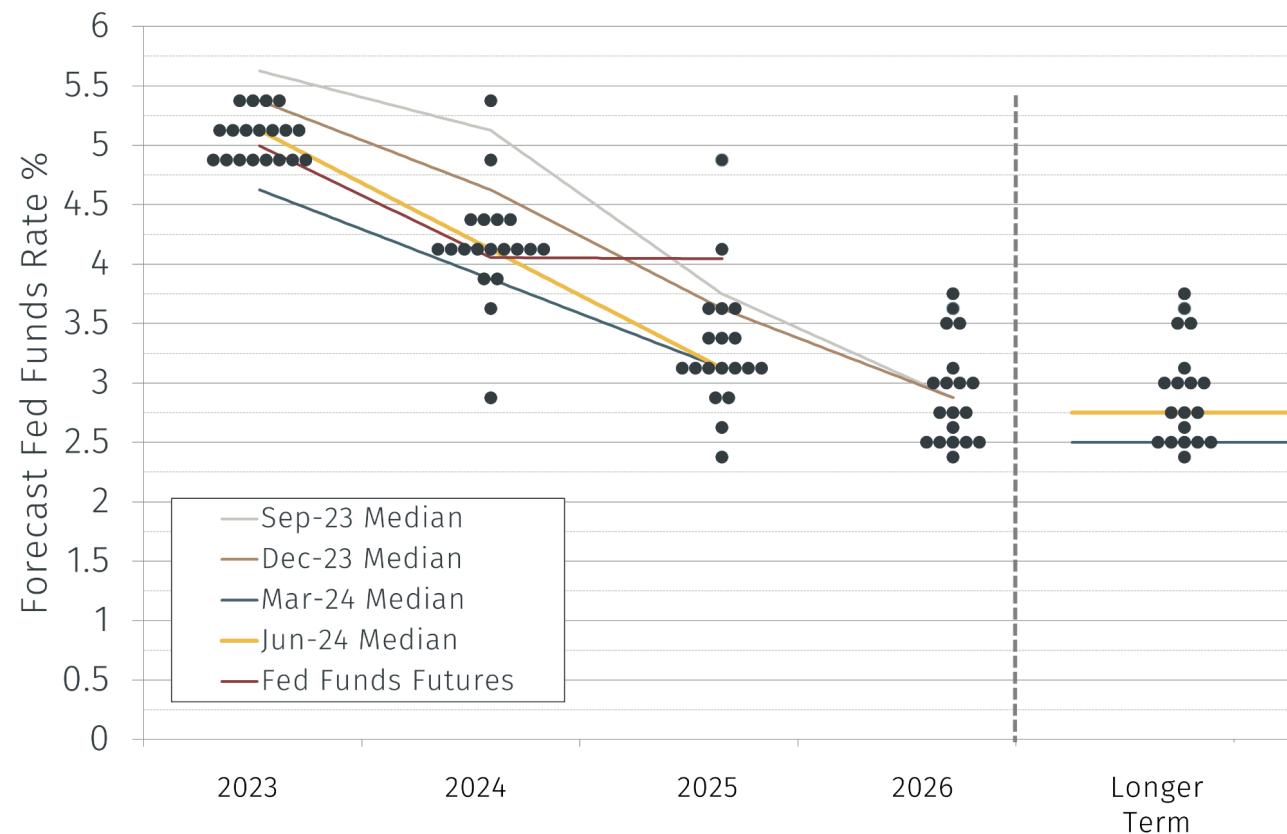


**SKKU  
Business School**

# Outline



# Prevalence of data or even big data in our lives



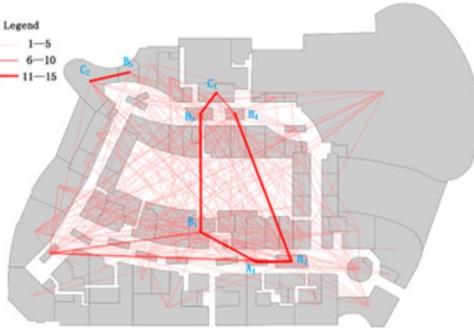
# Prevalence of data or even big data in our lives



# Prevalence of data or even big data in our lives



# Prevalence of data or even big data in our lives

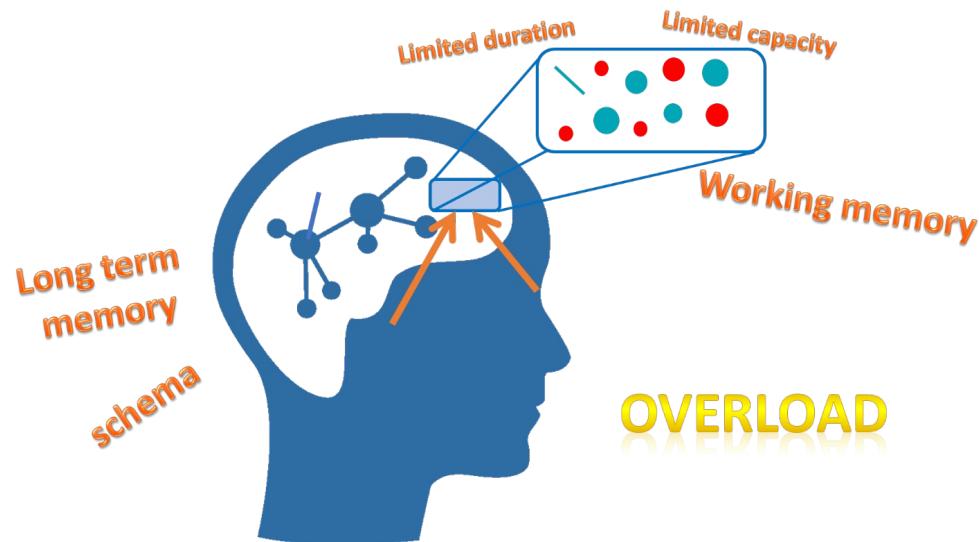


- Traditional data newly made available because of advancement in technology
- E.g., shopping trajectories
- New business contexts and interactions
- E.g., live streaming
- New data generated by AI trained using human data;
- New business contexts and interactions facilitated by AIGC
- New data by AI in the physical world
- E.g., self-driving cars, robots, drones

<https://www.youtube.com/watch?v=58BcWGojWHU#t=34m5s>



# Limited cognitive resources and attention



**Stereotyped bias** applies a widely held but fixed, oversimplified, and generalized image or idea of a particular type of person or thing.



**Similarity bias** leads people to connect with others who share similar interests, experiences, and backgrounds.



**Anchoring bias** causes us to rely too heavily on the first piece of information about a person or topic.



**Availability bias** distorts our view based on the most immediate information that readily comes to mind.



The **bandwagon effect** is groupthink or herd mentality, where people adopt behaviors and views of a dominant group.



**Confirmation bias** tends to interpret new evidence to validate one's existing beliefs.



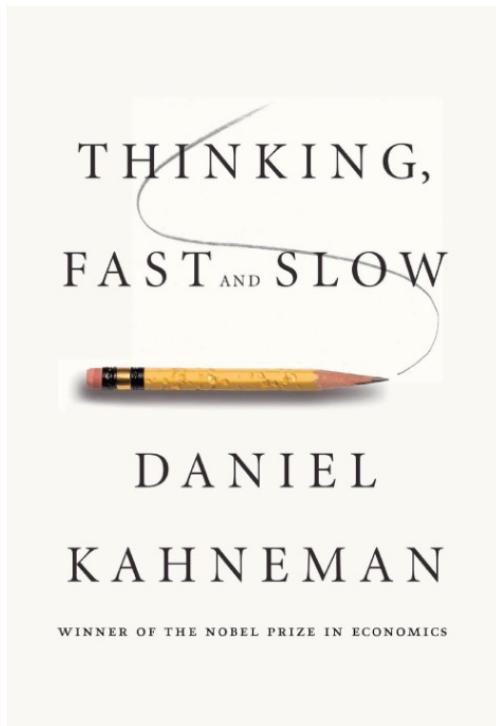
**Fundamental attribution error** tends to assume someone's actions are because of their character rather than circumstances.



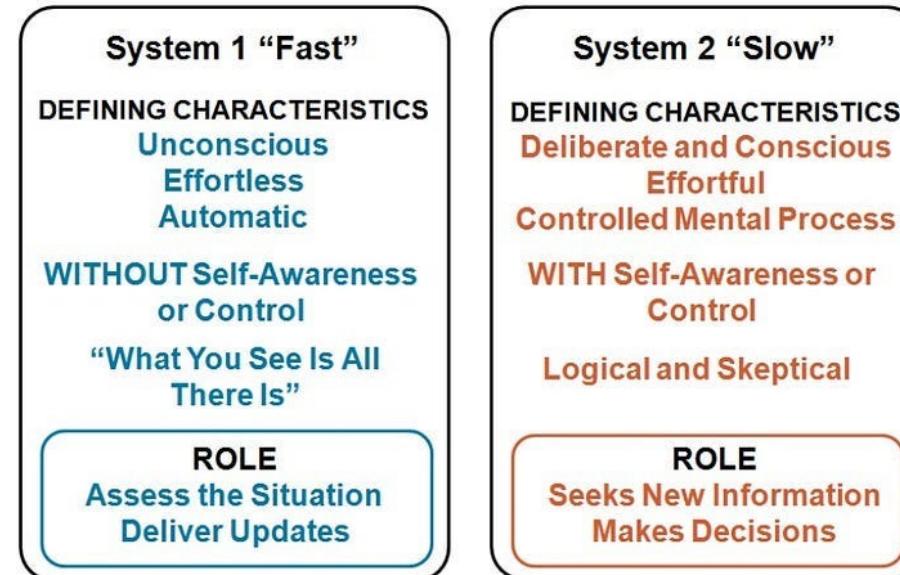
**Horns and halo effect** biases perceptions of someone because of opinions of that person's other traits.



# Limited cognitive resources and attention



## System 1 and 2 thinking



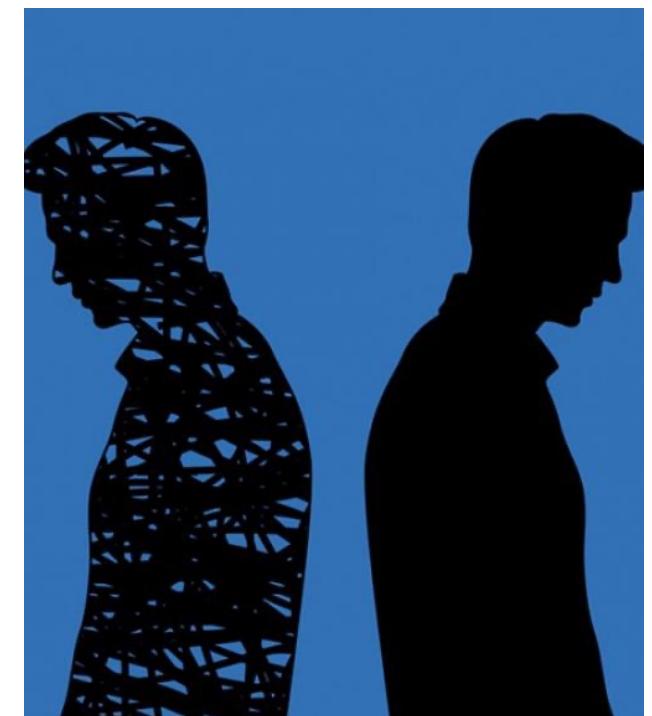
# Limited cognitive resources and attention



Attention/Eyeball Economy



## Urgency to make fast and ideally good decisions



Background	One Entity	Multi Entities	Relationships	Implications
<ul style="list-style-type: none"> <li>• Prevalence of data or even big data</li> <li>• Urgency to make fast and good decisions</li> <li>• Limited cognitive resource and attention</li> </ul>	<ul style="list-style-type: none"> <li>• Basic summary statistics</li> <li>• Bias in observability of estimates and forecasts</li> <li>• Bias in observed estimates and forecasts (e.g., fat tails)</li> </ul>	<ul style="list-style-type: none"> <li>• Dimension reduction</li> <li>• Market structure analysis with structured data</li> <li>• Market structure analysis with unstructured data</li> </ul>	<ul style="list-style-type: none"> <li>• Variable selection in one linear regression</li> <li>• Variable selection in multiple linear regressions</li> </ul>	<ul style="list-style-type: none"> <li>• Optimal aggregation rules</li> <li>• Impact of reviews on product sales</li> <li>• Impact of AIGS on review activities and content consumption</li> </ul>

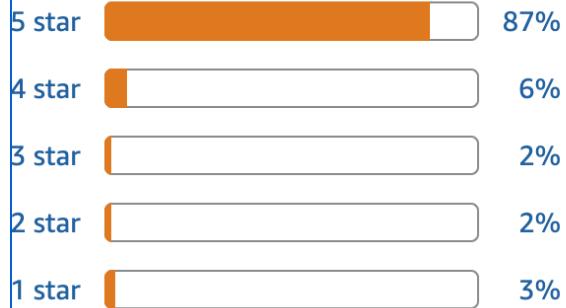


# Examples

## Customer reviews

★★★★★ 4.7 out of 5

146 global ratings



Angela

★★★★★ Owl Diaries ✨

Reviewed in the United States on February 6, 2024

Verified Purchase

My daughter loves these books. It has a series that's she's been reading since kindergarten on her own. I think it's a great way to start your child off with chapter books. I just love the funny names of the characters and the words they use. My daughter is in the first grade and she wants to read these over and over. 🎉

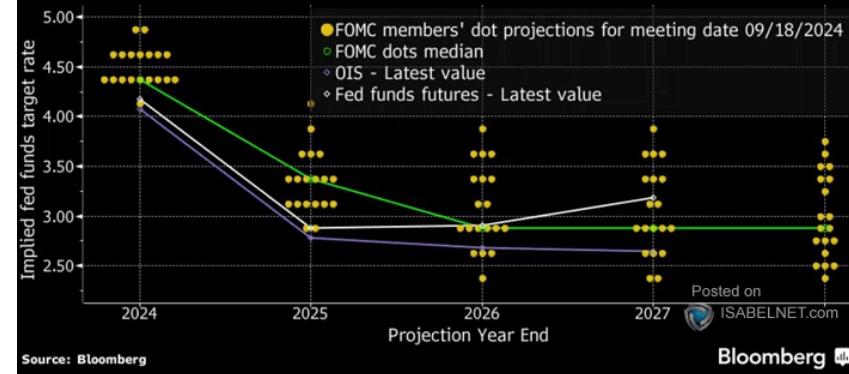
One person found this helpful

[Helpful](#) | [Report](#)

<https://www.amazon.com/Eva-President-Branches-Book-Diaries/dp/1338880276/>

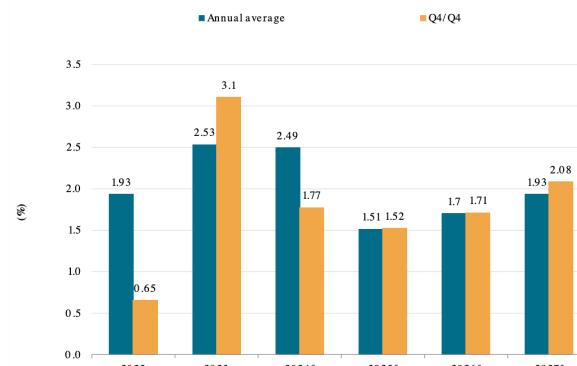


## The Fed's September Dot Plot



<https://www.isabelnet.com/implied-fed-funds-target-rate/>

Average U.S. real GDP growth is forecasted to reach 2.5% this year



<https://www.spglobal.com/ratings/en/research/articles/240326-economic-outlook-u-s-q2-2024-heading-for-an-encore-13048486>



THE HONG KONG  
POLYTECHNIC UNIVERSITY  
香港理工大學

# Basic summary statistics

- Number of observations
- Average of the observations
- Standard deviation of the observations
- *Correlation coefficient if observations of more than one variables*

Order statistics, oftentimes five-number summary

- Minimum of the observations
- The first quartile of the observations
- Median of the observations
- The third quartile of the observations
- Maximum of the observations

Example:

(14 73 28 42 9 85 31 19 67 51 22 98 11 46 75 39 88  
25 59 49 13 81 36 68 55 29 92 17 63 79)

## Summary Statistics

Statistic	Value
Number of Observations	30
Average	48.2
Standard Deviation	25.1
Minimum	9
Maximum	98
Median	46.5



# Not all estimates/ratings are observable

MARKETING SCIENCE  
 Vol. 31, No. 3, May–June 2012, pp. 372–386  
 ISSN 0732-2399 (print) | ISSN 1526-548X (online)



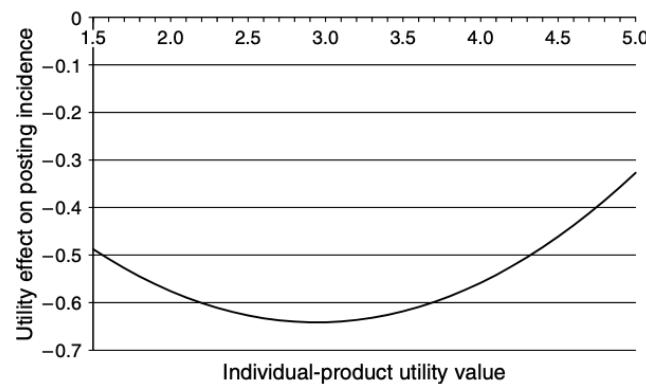
<http://dx.doi.org/10.1287/mksc.1110.0662>  
 © 2012 INFORMS

## Online Product Opinions: Incidence, Evaluation, and Evolution

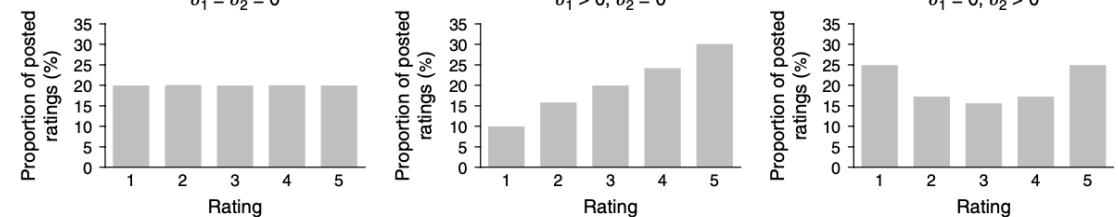
Wendy W. Moe  
 Robert H. Smith School of Business, University of Maryland, College Park, Maryland 20472,  
[wmoe@hsmith.umd.edu](mailto:wmoe@hsmith.umd.edu)

David A. Schweidel  
 Wisconsin School of Business, University of Wisconsin–Madison, Madison, Wisconsin 53706,  
[dschweidel@bus.wisc.edu](mailto:dschweidel@bus.wisc.edu)

**Figure 5** Role of Postexperience Evaluation in Rating Incidence



**Figure 2** Illustrative Distributions of Posted Ratings



# Estimates/ratings are oftentimes biased (1)

## Reviewing Experts' Restraint from Extremes and Its Impact on Service Providers

PETER NGUYEN  
XIN (SHANE) WANG  
XI LI  
JUNE COTTE

This research investigates reviewing experts on online review platforms. The main hypothesis is that greater expertise in generating reviews leads to greater restraint from extreme summary evaluations. The authors argue that greater experience generating reviews facilitates processing and elaboration and enhances the number of attributes implicitly contained in evaluations, which reduces the likelihood of assigning extreme summary ratings. The restraint-of-expertise hypothesis is tested across different review platforms (TripAdvisor, Qunar, and Yelp), shown for both assigned ratings and review text sentiment, and demonstrated both between (experts vs. novices) and within reviewers (expert vs. pre-expert). Two experiments replicate the main effect and provide support for the attribute-based explanation. Field studies demonstrate two major consequences of the restraint-of-expertise effect. (i) Reviewing experts (vs. novices), as a whole, have less impact on the aggregate valence metric, which is known to affect page-rank and conversion rates. (ii) Expert reviewers are more likely to provide service providers with their ratings. For service providers that generally provide mediocre (excellent) experiences, reviewing experts assign significantly higher (lower) ratings than novices. This research provides important caveats to the existing marketing practice of service providers incentivizing reviewing experts and provides strategic implications for how platforms should adopt rating scales and aggregate ratings.

**Keywords:** online word-of-mouth, expertise, user rating average, platform strategy, text analysis, sentiment analysis

Peter Nguyen ([p.nguyen@manchester.ac.uk](mailto:p.nguyen@manchester.ac.uk)) is an assistant professor of marketing at Miami University, Oxford, OH 45056, USA. Shane (Xin) Wang ([xwang@ivey.ca](mailto:xwang@ivey.ca)) is an associate professor of marketing and statistics, MBA '80 Faculty Fellow, at Ivey Business School, Western University, London, ON N6G 0W1, Canada. Xi Li ([xili44@cityu.edu.hk](mailto:xili44@cityu.edu.hk)) is an assistant professor of marketing at the City University of Hong Kong, Kowloon Tong, Hong Kong, China. June Cotte ([jcotte@ivey.ca](mailto:jcotte@ivey.ca)) is the professor of marketing, Scott & Melissa Beattie Professorship in Marketing, at Ivey Business School, Western University, London, ON N6G 0W1, Canada. This research was conducted prior to Peter's death. This article is based on the first author's dissertation. The authors thank seminar participants at Duke University, Queen's University, and the University of Connecticut, as well as participants in the AMAS-CBSIG conference in Bern, for their valuable feedback on earlier versions of this article. This work was supported by the Social Sciences and Humanities Research Council (435-2019-0597). **Supplementary materials** are included in the web appendix accompanying the online version of this article.

**Editor:** J. Jeffrey Inman

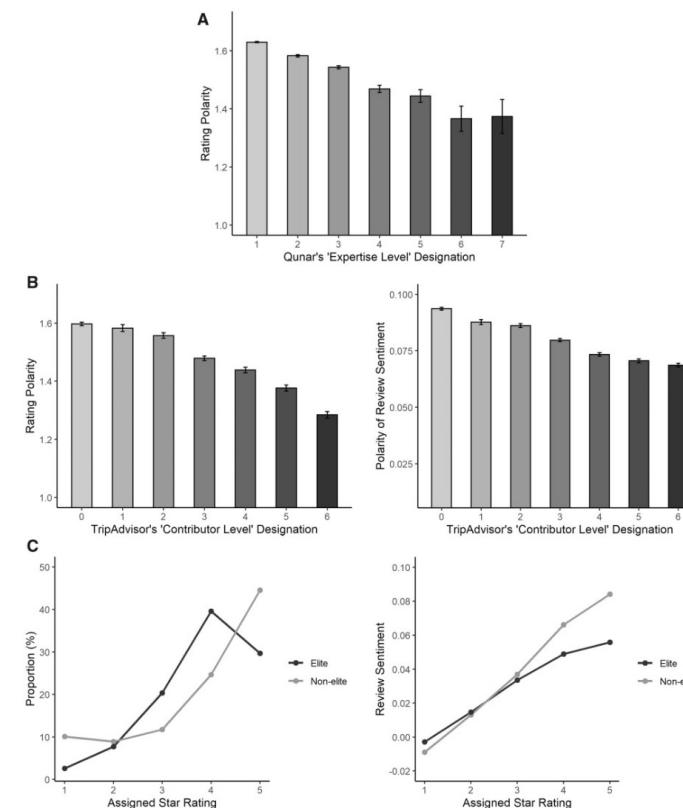
**Associate Editor:** Andrew T. Stephen

**Advance Access publication 15 July 2020**

Consumers rely on the opinions and recommendations of others. Many of these recommendations have come from expert professionals (e.g., sommeliers, movie critics). Over the past couple of decades, the world has seen the rise of online reviews, where consumers not only rely on other consumers' experiences but also share their own. Online review platforms now recognize their top users as reviewing "experts." For example, Yelp has its "Elite" status, TripAdvisor has its "Contributor Level," Google has its "Local Guide" badges, and Amazon has its "Amazon Vine Program." Given that consumers are increasingly both sharing and consuming reviews, understanding the nature of reviewing experts has become an important topic in consumer research.

The study of online reviewing experts is particularly important for *service providers*, such as hotels and restaurants. Many businesses incentivize, by quite literally

FIGURE 1  
POLARITY OF EVALUATIONS AS A FUNCTION OF PLATFORM-DEFINED REVIEWING EXPERTISE. (A) QUNAR (STUDY 1). (B) TRIPADVISOR (STUDY 3). (C) YELP (STUDY 4) ON



# Estimates/ratings are oftentimes biased (2)

Fat Tails in Human Judgment:

Empirical Evidence and Implications for  
the Aggregation of Estimates and Forecasts

Miguel Sousa Lobo\*

Dai Yao†

September 18, 2024

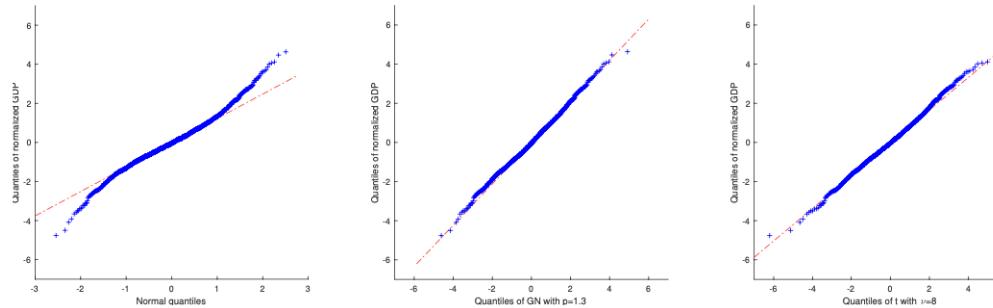


Figure 1: Quantile-quantile plots of USA gross domestic product growth forecasts by economists (after each quarter's forecasts were normalized to have zero mean and unit variance), against a normal distribution, against a generalized normal distribution with shape parameter  $p = 1.3$ , and against a Student's  $t$ -distribution with the number of degrees of freedom  $\nu = 8$ .

$$f(x) = \frac{1}{2s\Gamma(1+1/p)} \exp\left\{-\left|\frac{x-u}{s}\right|^p\right\}.$$

Data set	$N$	$n$	Skew.	Kurt.	$p_{kurt}$	$k$	$p_{Bayes}$	$p_{low}$	$p_{high}$
UNHI	2017	18	-0.91	5.3	1.09	3	0.95	0.86	1.04
UNLO	2022	18	-0.72	4.1	1.36	3	1.52	1.41	1.64
FXHI	2003	18	0.38	3.5	1.64	39	1.69	1.54	1.85
FXLO	2025	18	-0.27	4.1	1.36	39	1.35	1.22	1.49
SCALE	3618	34	0.42	3.9	1.43	71	1.66	1.55	1.77

Data set	$N$	$m$	$n$	Homogeneous judges				Heterogeneous judges			
				$k$	$p_{Bayes}$	$p_{low}$	$p_{high}$	$k$	$p_{Bayes}$	$p_{low}$	$p_{high}$
GDP	3682	125	70	141	1.23	1.15	1.31	395	1.81	1.65	1.96
GDPa	2158	99	42	85	1.12	1.03	1.21	287	1.67	1.48	1.89
NFARM	8783	127	170	341	1.17	1.12	1.22	599	1.64	1.55	1.74

Table 7: Estimates of the GN shape parameter  $p$  for economic forecasts, controlling for heterogeneity across forecasters.

Table 6: Estimates of the GN shape parameter  $p$  for trivia questions.



\*miguel.lobo@insead.edu, Associate Professor of Decision Sciences at INSEAD.

†dai@yaod.ai, Associate Professor at The Hong Kong Polytechnic University.

# How about summaries of data in unstructured form



Gheen

★★★★★ **Beauty -- In Every Way**

Reviewed in the United States on November 18, 2024

Size: 42 Inch | Style: TV Only | **Verified Purchase**

Buying a TV for any room has few considerations - but all important. They are the appropriate size of the TV to its service space, the quality of color production, the quality of sound, ease and efficiency of controls and the ability to set install and setup the TV without having to employ and engineer. This LG Class OLED exudes excellence in all areas at a reasonable cost. Wall mount installation and operational setup took 15 minutes each. The color quality is stunningly beautiful, and the program walks you through making fine adjustments to meet your color tastes. The sound quality was excellent, and so beautiful that a plan to add a sound bar were scrapped. The TV controller is easy to manage and instantly responsive. At age 74 and having owned a plethora of TVs across the years, the LG is the best and most beautiful buy.

11 people found this helpful

Helpful

| Report



Tony S

★★★★★ **Ultimate Gaming Experience – Stunning Visuals and Smart Features!**

Reviewed in the United States on October 22, 2024

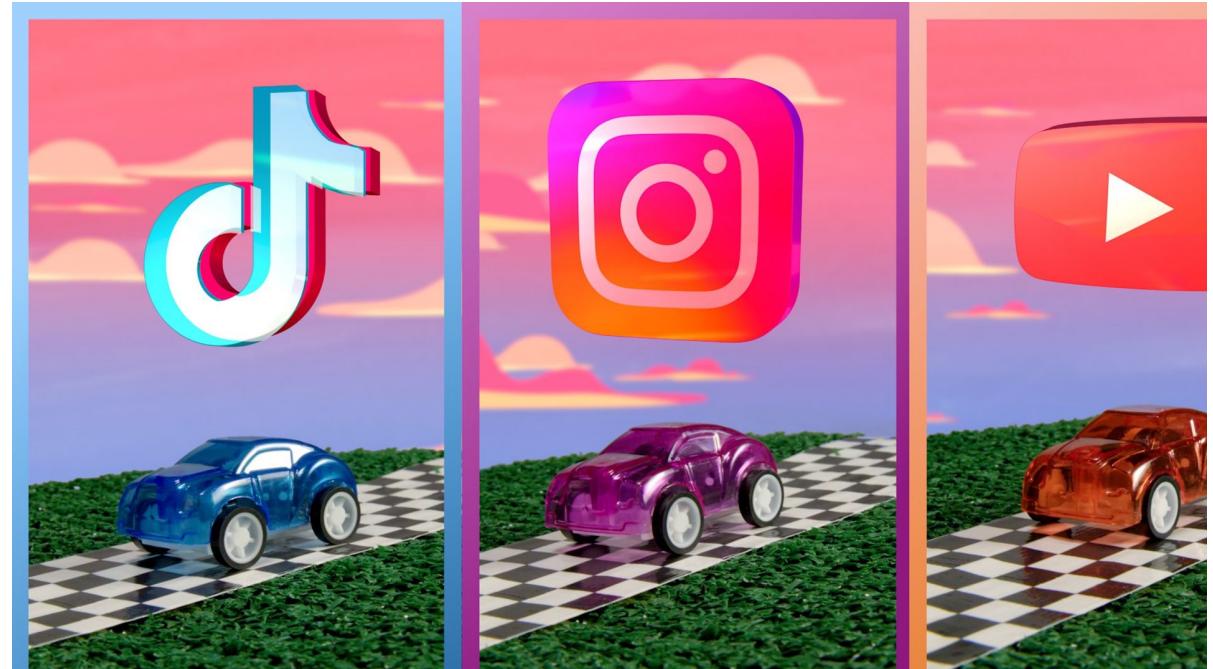
Size: 48 Inch | Style: TV Only | **Verified Purchase**

The LG 48-Inch OLED evo C4 is hands down the best TV I've used for gaming. The picture quality is incredible, with deep blacks, vibrant colors, and smooth motion that make every game look stunning. The 4K processor ensures everything runs seamlessly, and the OLED technology eliminates motion blur, making it perfect for fast-paced games.

What really impressed me is the low input lag and the built-in gaming features like NVIDIA G-Sync and FreeSync support. It provides a responsive experience, giving me an edge in competitive play. The size is also ideal for immersive gaming without being overwhelming.

The AI-powered Magic Remote and Alexa integration make it incredibly convenient to control both the TV and my smart home devices. Switching between gaming and streaming apps is quick, and the user interface is intuitive.

For the price, this TV offers exceptional value. Whether you're a casual or serious gamer, the LG OLED evo C4 delivers a next-level experience. Highly recommend it!



Texts



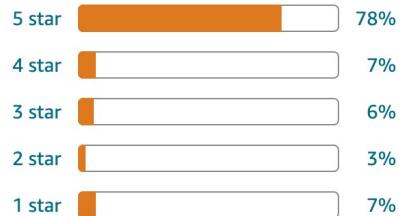
Short and long videos

# Amazon's solution: AI-Generated Summaries (AIGS)

## Customer reviews

★★★★★ 4.5 out of 5

261 global ratings



⌄ How customer reviews and ratings work

## Review this product

Share your thoughts with other customers

Write a customer review

## Customers say

Customers like the picture quality, ease of setup, and sound quality of the television. They mention it's divine, vibrant, and 4K content looks stunning. Some appreciate the value for money, saying it gives great bang for the buck. Customers also appreciate the build quality and speed.

AI-generated from the text of customer reviews

### Select to learn more

- ✓ Picture quality | ✓ Ease of setup | ✓ Sound quality | ✓ Value for money | ✓ Brightness |
- ✓ Build quality | ✓ Speed

## Reviews with images

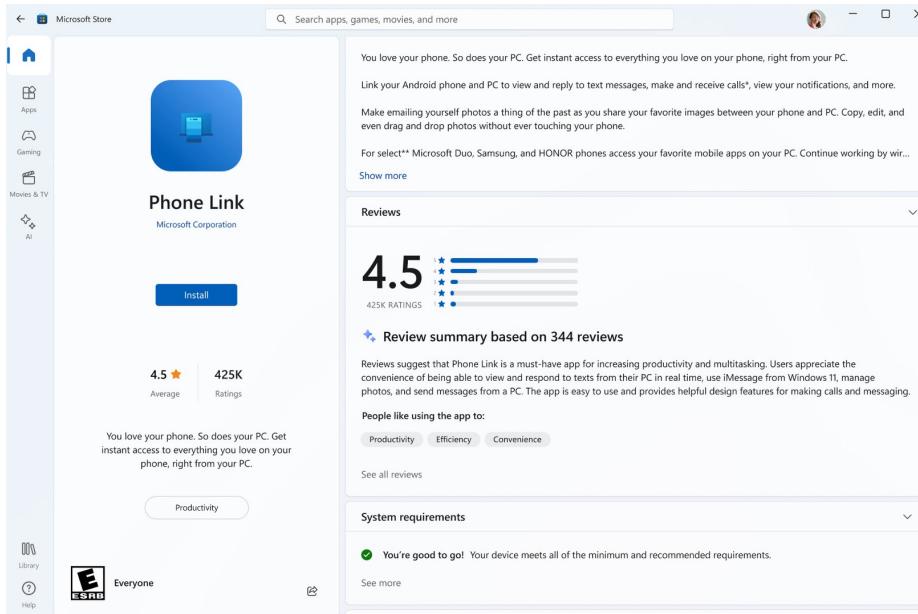
See all photos >



<https://www.amazon.com/LG-55-Inch-Processor-AI-Powered-OLED55C4PUA/dp/B0CVRDK4P6/>



# Microsoft Store, Newegg are also embracing the trend



<https://www.theverge.com/2023/5/23/23732821/microsoft-store-bing-ai-hub-apps-build>

A screenshot of a Newegg product page for the 'Intel Core i9-13900K - Core i9 13th Gen Raptor Lake 24-Core (8P+16E) P-core Base Frequency: 3.0 GHz E-core Base Frequency: 2.2 GHz'. The page includes a navigation bar with 'Overview', 'Specs', 'Reviews', 'Q &amp; A', 'Compare Products', 'Warranty &amp; Returns', and 'More Buying Options'. Below the navigation is a section titled 'WHAT REVIEWERS ARE SAYING ABOUT THE PRODUCT' with a summary: 'The Intel Core i9-13900K is highly praised for its fast performance, compatibility with existing Z690 chipsets, and improved memory controller. It is considered a great upgrade from older CPUs and is particularly well-suited for gaming and productivity tasks. However, the CPU is known to run hot and consume a lot of power, requiring a good cooling solution. Some users have also mentioned limited cooler compatibility and high price as potential disadvantages. Overall, the i9-13900K is a powerful CPU that delivers excellent performance, but it may require careful consideration of cooling and power requirements.' There are also sections for 'Pros' (Fast performance, Good for gaming, Good for overclocking, Efficient performance) and 'Cons' (Runs hot). At the bottom right are 'FEEDBACK', 'Like', and 'Share' buttons.

<https://www.theverge.com/2023/8/9/23825805/newegg-chatgpt-review-summaries-ai>



# Video platforms are also embracing it

小米 14 体验，影像性能双进化

410 1 2023-10-27 15:50:11 未经作者授权，禁止转载



IT之家 发消息  
爱科技，爱这里 - 前沿科技人气平台。合作联系：v@r...  
+ 充电 + 关注 5.7万

已为你生成视频总结 ①

IT之家拿到的黑色小米14手机的外观设计、摄像头模组、处理器、内存、电池和软件优化等方面。小米14采用了亮面不锈钢中框和磨砂玻璃四曲面后盖,摄像头模组采用圆角矩形设计,与徕卡深度合作。处理器方面,小米14搭载了骁龙八进三和小米自家澎湃OS的首发,硬件配置拉满,软件优化十分丝滑。此外,视频还提到了小米14的游戏性能、生态优势等方面。

- 小米14手机的外观设计、摄像头、处理器、软件优化等方面的详细介绍。

00:01 小米14发布,采用黑色亮面不锈钢中框和磨砂玻璃四曲面后盖  
01:08 小米14镜头采用徕卡不同镜头的标识,提升影像素质  
02:49 骁龙八进三处理器和小米自家澎湃OS的首发,硬件配置拉满,软件优化良好



<https://www.ithome.com/0/728/123.htm>

# A lot of tools available within/beyond main platforms

## YouTube Video Summarizer

Get YouTube transcript and use AI to summarize YouTube videos in one click for free online with NoteGPT's YouTube summary tool.

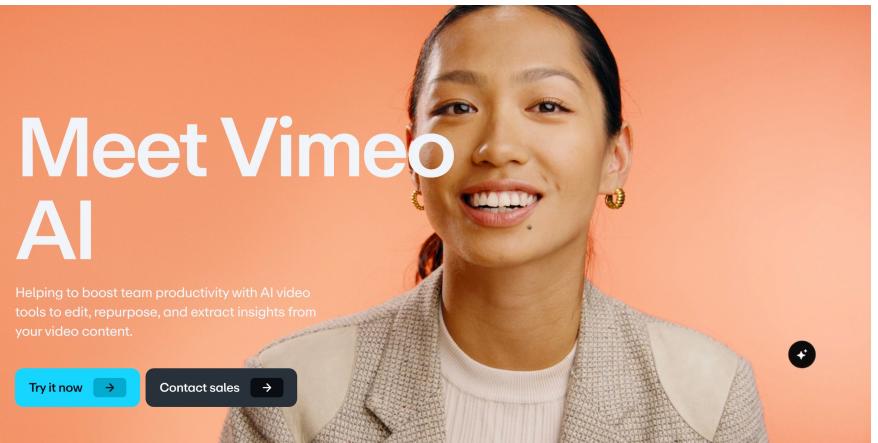
Want to summarize local Video/Audio files? [Go to Workspace to summarize](#)

## YouTube Summary with ChatGPT Online and Free

Summarize YouTube videos in seconds

Paste a YouTube video url here

Summarize



The landing page for Vimeo AI features a large, smiling woman of Asian descent on the right side. To her left, the text "Meet Vimeo AI" is displayed in large, white, sans-serif font. Below this, a smaller text block reads: "Helping to boost team productivity with AI video tools to edit, repurpose, and extract insights from your video content." At the bottom, there are two buttons: a blue "Try it now" button and a black "Contact sales" button.

## Vimeo Summarizer

Summarize Vimeo videos in seconds with ScreenApp's AI

Paste a URL (YouTube, Vimeo, Facebook, TikTok)

Import from URL



Loved by over 1 million users

# Summary of “live” content: Zoom



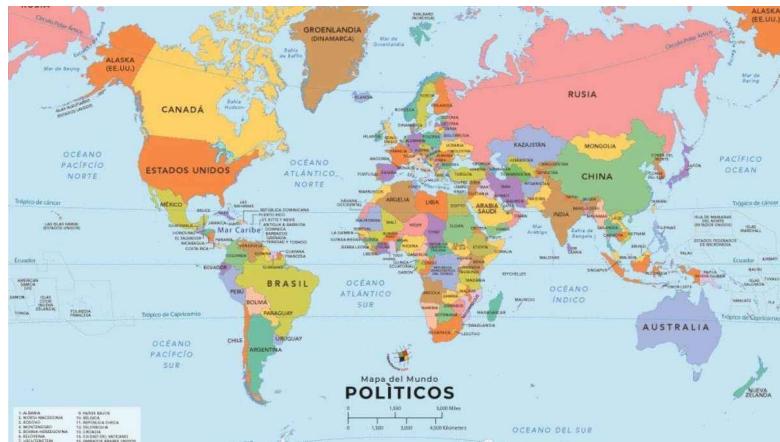
<https://www.zoom.com/en/blog/zoom-ai-companion-getting-started-guide/>



Background	One Entity	Multi Entities	Relationships	Implications
<ul style="list-style-type: none"> <li>• Prevalence of data or even big data</li> <li>• Urgency to make fast and good decisions</li> <li>• Limited cognitive resource and attention</li> </ul>	<ul style="list-style-type: none"> <li>• Basic summary statistics</li> <li>• Bias in observability of estimates and forecasts</li> <li>• Bias in observed estimates and forecasts (e.g., fat tails)</li> </ul>	<ul style="list-style-type: none"> <li>• Dimension reduction</li> <li>• <b>Market structure analysis with structured data</b></li> <li>• <b>Market structure analysis with unstructured data</b></li> </ul>	<ul style="list-style-type: none"> <li>• Variable selection in one linear regression</li> <li>• Variable selection in multiple linear regressions</li> </ul>	<ul style="list-style-type: none"> <li>• Optimal aggregation rules</li> <li>• Impact of reviews on product sales</li> <li>• Impact of AIGS on review activities and content consumption</li> </ul>



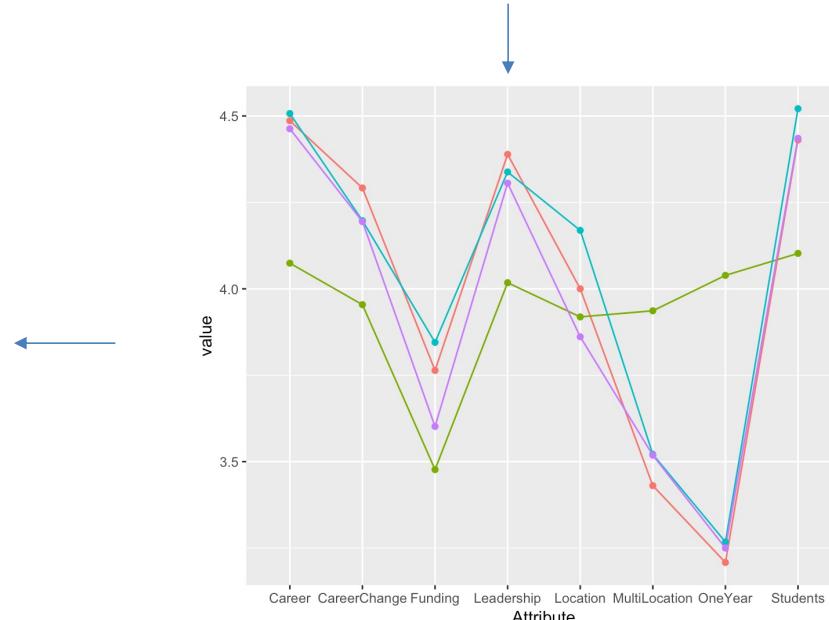
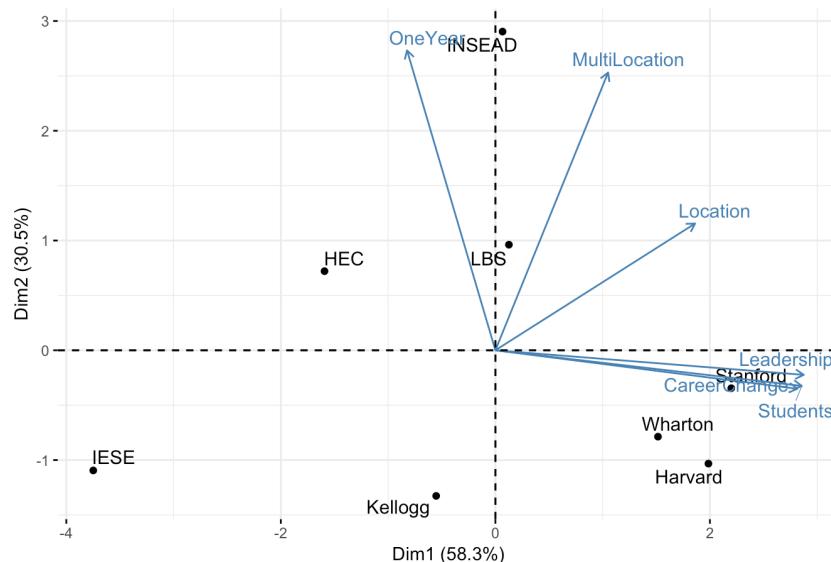
# Examples



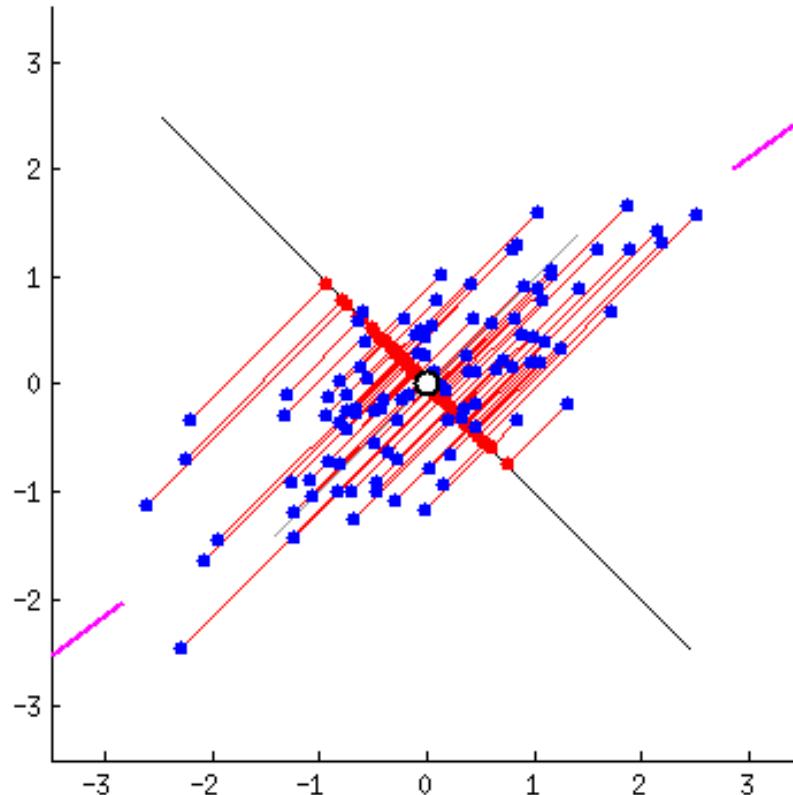
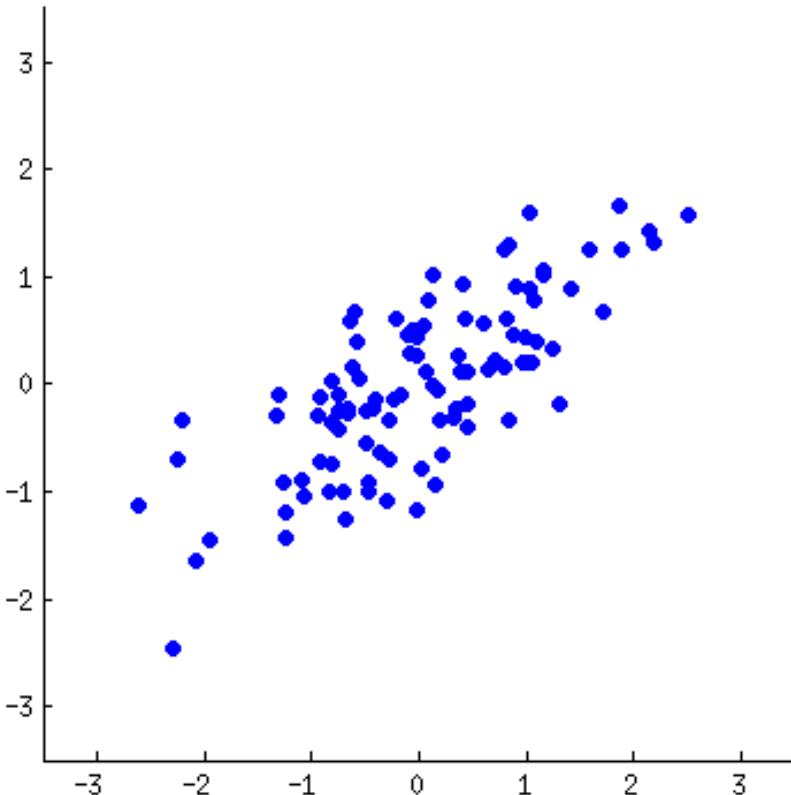
# Aggregation of information about multiple entities

SchoolRated	Q17_1	Q17_2	Q17_3	Q17_4	Q17_5	Q17_6	Q17_7	Q17_8	Q17_9	Q17_10	Q17_11	Q17_12	Q17_13
8	3	5	5	4	4	4	2	2	3	4	5	3	5
7	4	5	5	4	4	2	2	2	3	4	4	3	5
1	2	5	5	5	5	5	4	2	5	5	5	3	5
4	5	5	5	5	3	5	3	5	5	4	5	5	5
3	3	3	3	3	3	3	3	3	3	3	3	3	3
4	3	3	3	3	3	3	3	3	3	3	3	3	3
6	3	5	5	5	5	5	3	3	3	5	4	3	4
8	3	5	5	5	5	3	5	3	3	5	4	3	4
4	3	4	4	5	3	5	3	3	4	4	3	4	4

Attribute	Harvard	HEC	IESE	INSEAD	Kellogg	LBS	Stanford	Wharton
OneYear	3.208333	3.605263	3.40	4.038869	3.354839	3.679487	3.267606	3.250000
Students	4.430556	3.745614	3.58	4.102473	4.064516	4.102564	4.521127	4.435185
Leadership	4.388889	3.807018	3.42	4.017668	4.010753	3.974359	4.338028	4.305556
CareerChange	4.291667	3.605263	3.39	3.954064	3.978495	3.858974	4.197183	4.194444
Location	4.000000	3.973684	3.55	3.918728	3.709677	4.269231	4.169014	3.861111
MultiLocation	3.430556	3.508772	3.23	3.936396	3.258065	3.461538	3.521127	3.518519



# Principal component analysis



which is very similar to factor analysis, though the two are fundamentally different.

<https://www.theanalysisfactor.com/the-fundamental-difference-between-principal-component-analysis-and-factor-analysis/>



# Principal component analysis: An illustration

Who's #1? INSEAD, Harvard, Wharton, LBS?

## ***Who's who?***

"The Business School for the World"

"Impact the way the world does business"

"Global focus and learning experience throughout programs"

"We educate leaders who make a difference in the world"

"Our ultimate goal: to create a global world that has meaning for us all"

"The Global Business School"

## ***Who's who?***

**INSEAD**  
The Business School  
for the World®

**London**  
Business  
School

 THE UNIVERSITY OF MELBOURNE | MELBOURNE BUSINESS SCHOOL

 HARVARD  
BUSINESS SCHOOL

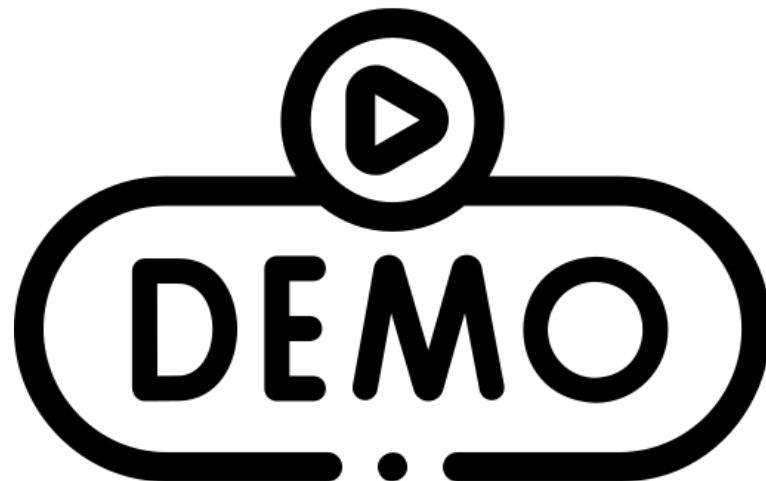
**ESSEC**  
BUSINESS SCHOOL

 HULT  
INTERNATIONAL  
BUSINESS SCHOOL



 THE HONG KONG  
POLYTECHNIC UNIVERSITY  
香港理工大學

# Principal component analysis: An illustration



# Two main approaches of market structure analysis

Brands as collections of *disaggregate* attributes, using principal component analysis

Basic data structure

Attribute	Harvard	HEC	IESE	INSEAD	Kellogg	LBS	Stanford	Wharton
OneYear	3.20833	3.60526	3.48	4.03886	3.35483	3.67948	3.267606	3.250000
Students	4.43055	3.74561	3.58	4.102473	4.064516	4.102564	4.521127	4.435185
Leadership	4.38889	3.807018	3.42	4.017668	4.018751	3.974359	4.338028	4.305556
CareerChange	4.291667	3.605263	3.39	3.954664	3.978495	3.858974	4.197183	4.194444
Location	4.000000	3.973684	3.55	3.918728	3.709677	4.269231	4.169014	3.861111
Multilocation	3.430556	3.508772	3.23	3.936396	3.258065	3.461538	3.521127	3.518519



John, Loken, Kim,  
Monga (JMR 2006)



Culotta and Cutler  
(MKSC 2016)



Liu, Dzyabura, Mizik  
(MKSC 2020)  
Dzyabura, Peres (JM 2021)



Li, Castelo, Katona,  
Sarvary (MKSC 2024)

Brands as *aggregate* objects, using clustering or multidimensional scaling based on similarity matrix



Kim, Albuquerque,  
Bronnenberg (JMR 2011)  
Ringel, Skiera (MKSC 2016)



Gabel, Guhl, Klapper  
(JMR 2019)  
Ruis, Athey, Blei (AAS  
2020)



Chintagunta (MKSC 1998)  
France and Ghose (MKSC  
2016)



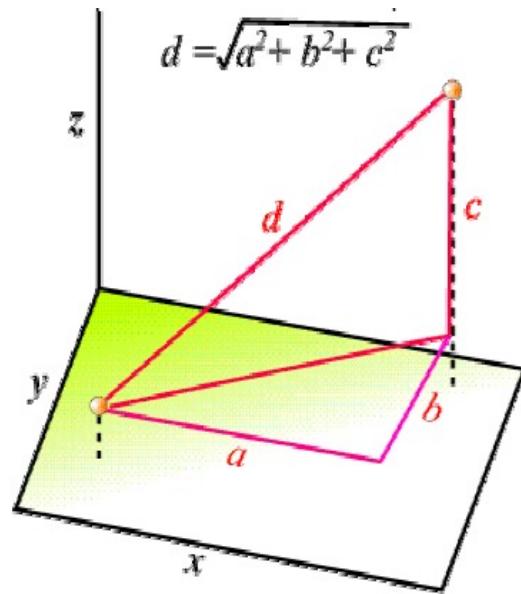
Lee and Bradlow (JMR 2011)  
Netzer, Feldman, Goldenberg,  
Fresko (MKSC 2012)  
Tirunillai and Tellis (JMR 2014)  
Matthe, Ringel, Skiera (MKSC 2022)  
Zhang, Kim, Xing (KDD 2015)



Yang, Zhang, Kannan  
(JM 2021)

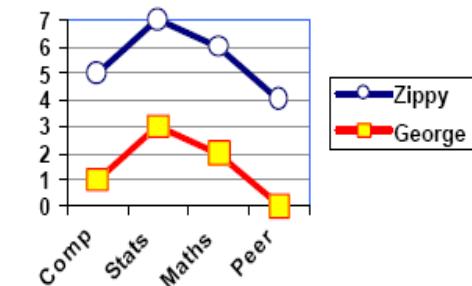
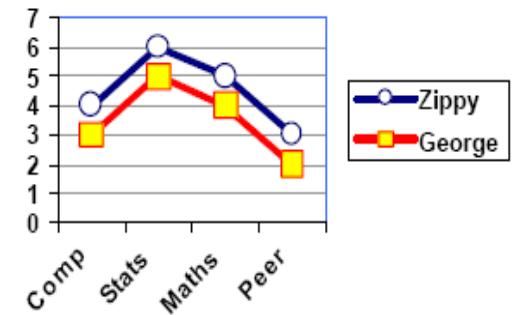


# Clustering – Similarity measures



$$\mathbf{D} = \begin{bmatrix} 0 & d_{12} & d_{13} \\ d_{21} & 0 & d_{23} \\ d_{31} & d_{32} & 0 \end{bmatrix}$$

$$\mathbf{r} = \begin{bmatrix} 1 & r_{12} & r_{13} \\ r_{21} & 1 & r_{23} \\ r_{31} & r_{32} & 1 \end{bmatrix}$$

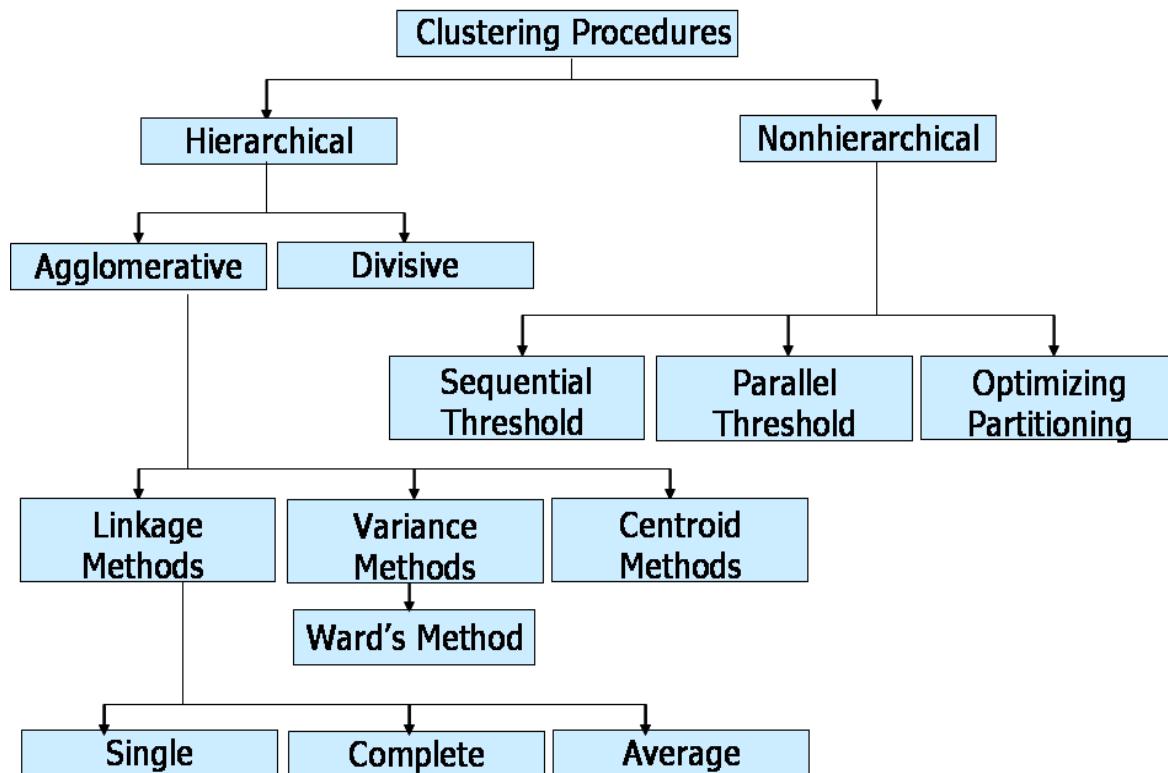


- Euclidian distance, city-block or Manhattan distance, Chebyshev distance
- Variable standardization before calibration

- Correlation is simple to calculate
- The elevation of score might be an issue



# Clustering – Methods



**Hierarchical clustering:** is characterized by the development of a hierarchy or tree-like structure. Hierarchical methods can be agglomerative or divisive.

**Nonhierarchical clustering:** These methods are frequently referred to as  $k$ -means clustering. These methods include sequential threshold, parallel threshold, and optimizing partitioning.

**Agglomerative (hierarchical) clustering:** starts with each object in a separate cluster. Clusters are formed by grouping objects into bigger and bigger clusters. This process is continued until all objects are members of a single cluster.

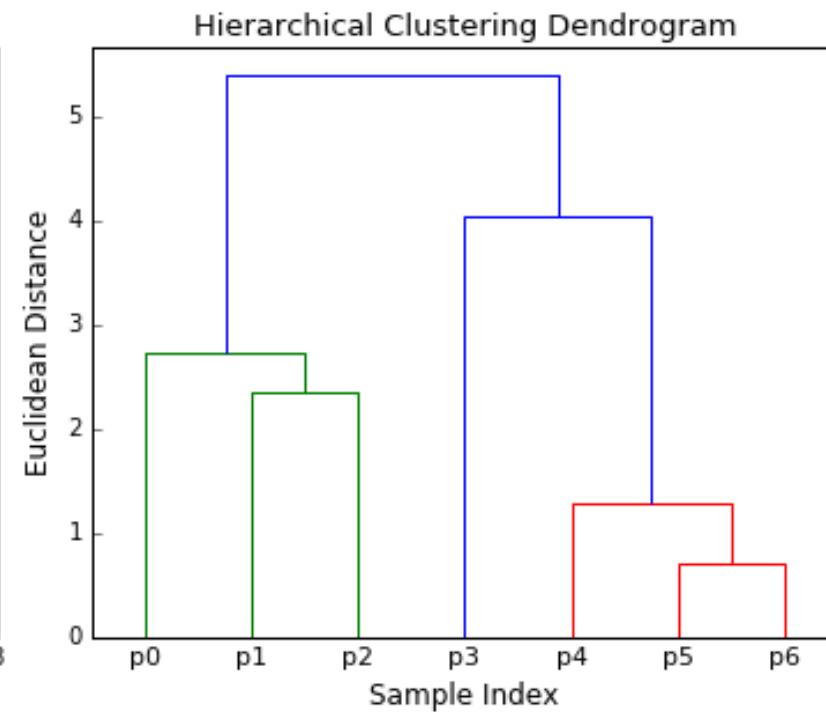
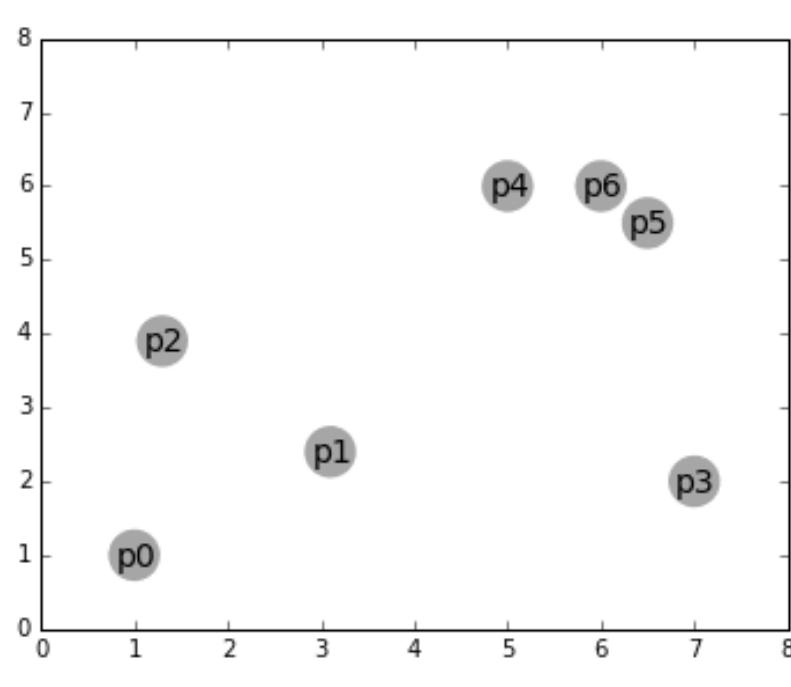
**Divisive (hierarchical) clustering:** starts with all the objects grouped in a single cluster. Clusters are divided or split until each object is in a separate cluster.

Agglomerative methods are commonly used in marketing research. They consist of

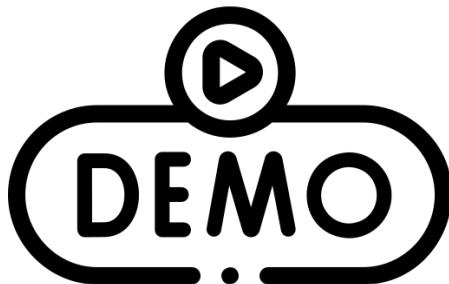
- linkage methods,
- variance methods or error sums of squares, and
- centroid methods.



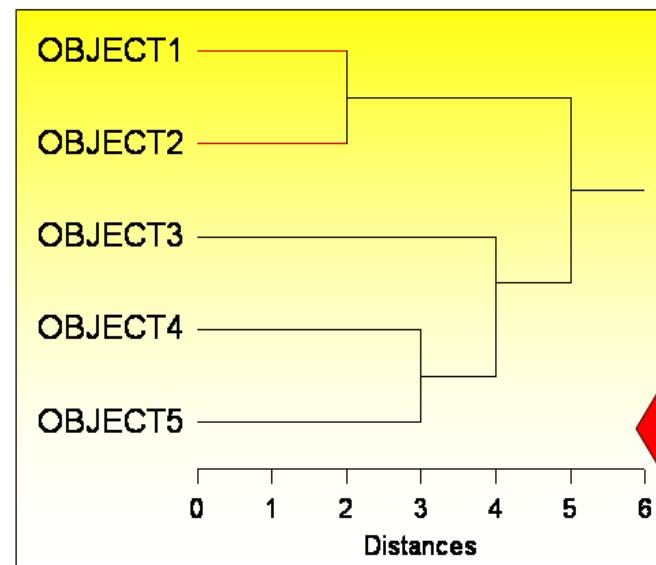
## Clustering – hierarchical



## Clustering – hierarchical: An illustration



### Simple joining (nearest neighbour)



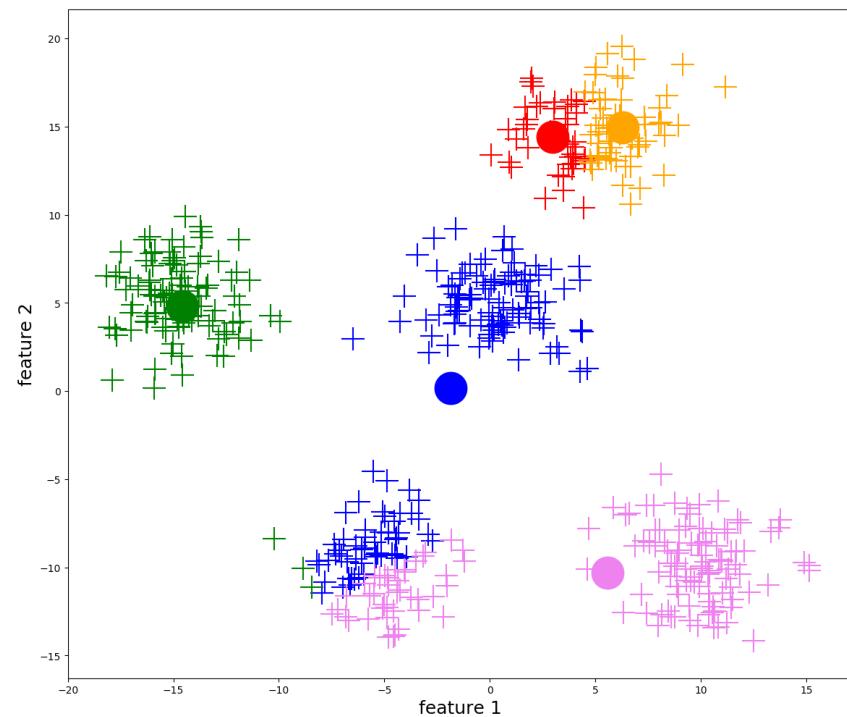
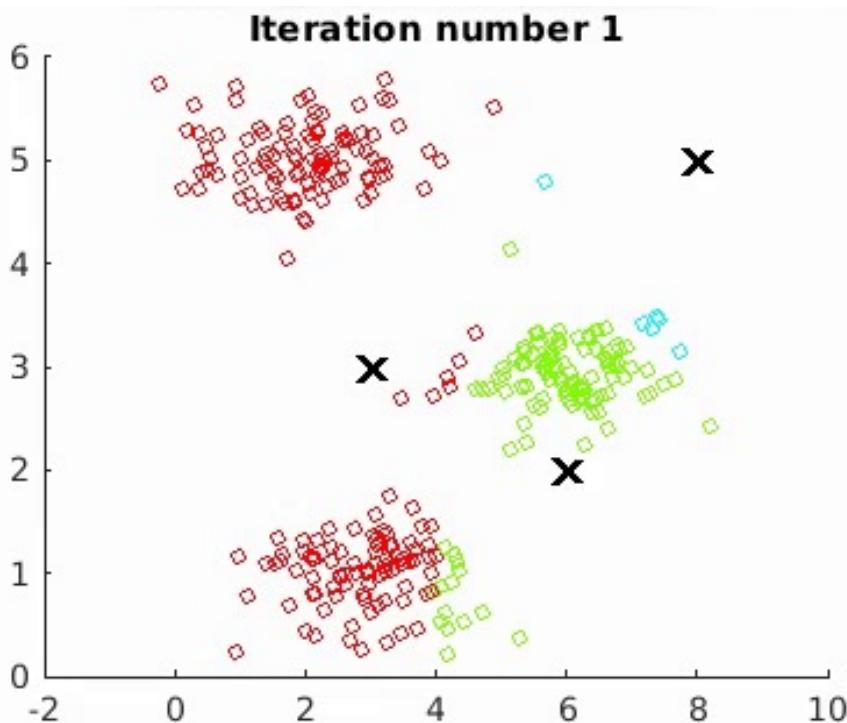
Object	1	2	3	4	5
1					
2		2			
3			6	5	
4				10	9
5					3

**Distance matrix**

Distance	Cluster
0	1, 2, 3, 4, 5
2	(1, 2), 3, 4, 5
3	(1, 2), 3, (4, 5)
4	(1, 2), (3, 4, 5)
5	(1, 2, 3, 4, 5)



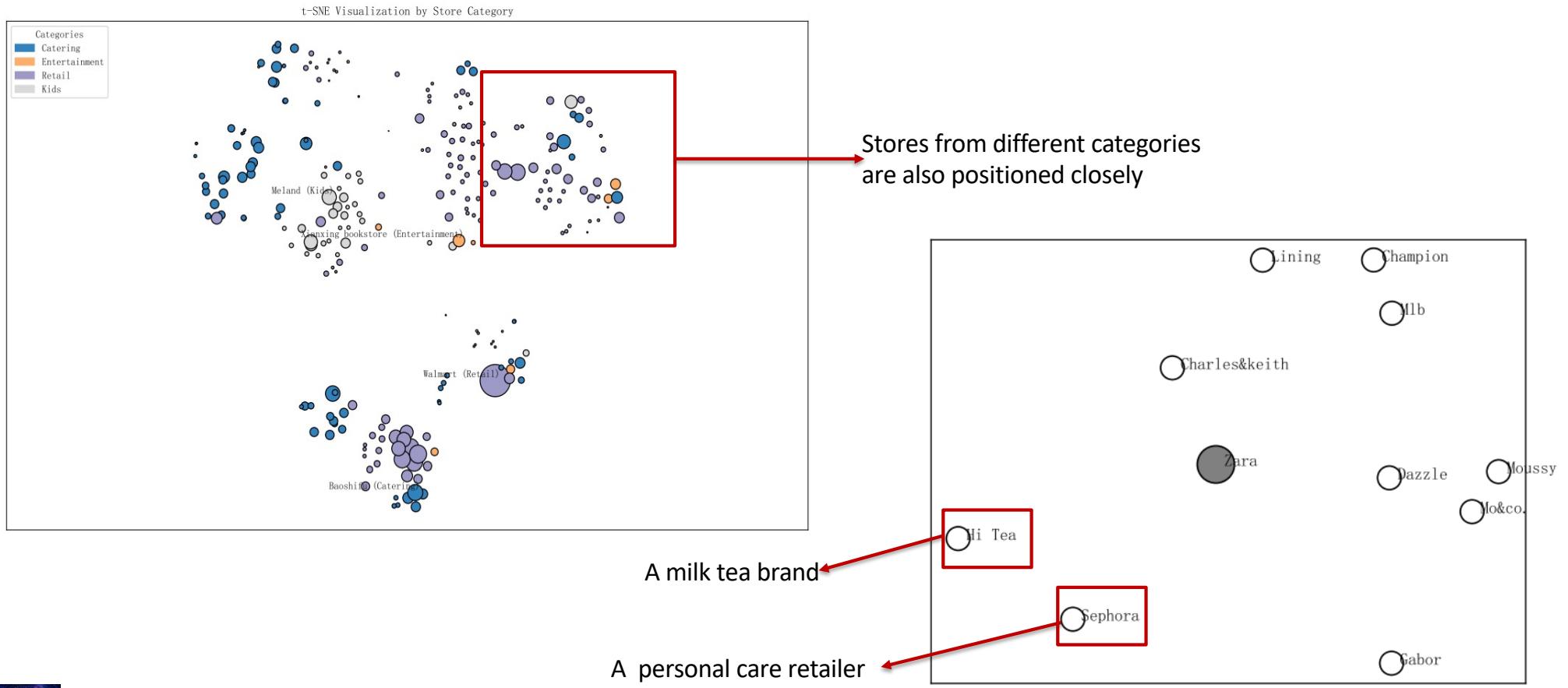
## Clustering – non-hierarchical, e.g., kmeans



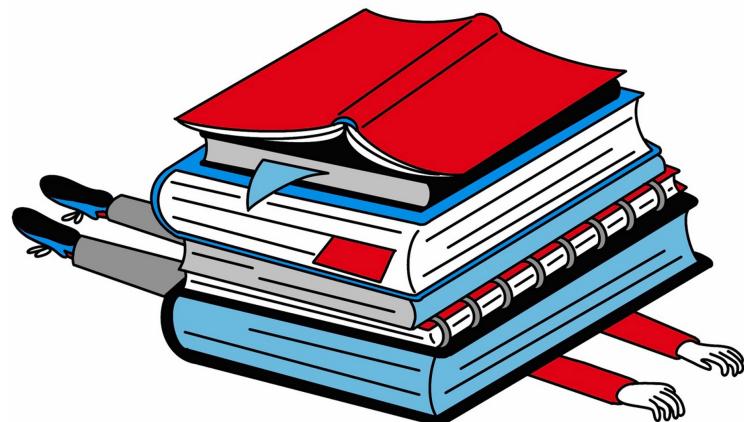
# Digitalization of Retailing Spaces and Shopping Trajectories



# t-SNE visualization of the shopping mall market structure



# Logistics before next class (3 Dec 2024)



Very important!

Submit your answers by 11:59pm, 2 Dec 2024

Think about the implications (we'll discuss some in the last section in next class)



Background	One Entity	Multi Entities	Relationships	Implications
<ul style="list-style-type: none"> <li>• Prevalence of data or even big data</li> <li>• Urgency to make fast and good decisions</li> <li>• Limited cognitive resource and attention</li> </ul>	<ul style="list-style-type: none"> <li>• Basic summary statistics</li> <li>• Bias in observability of estimates and forecasts</li> <li>• Bias in observed estimates and forecasts (e.g., fat tails)</li> </ul>	<ul style="list-style-type: none"> <li>• Dimension reduction</li> <li>• Market structure analysis with structured data</li> <li>• Market structure analysis with unstructured data</li> </ul>	<ul style="list-style-type: none"> <li>• Variable selection in one linear regression</li> <li>• Variable selection in multiple linear regressions</li> </ul>	<ul style="list-style-type: none"> <li>• Optimal aggregation rules</li> <li>• Impact of reviews on product sales</li> <li>• Impact of AIGS on review activities and content consumption</li> </ul>



# Aggregation of information about some relationships

$$Y_i = \beta_0 + \beta_1 X_i \quad \xrightarrow{\hspace{1cm}} \quad \arg \min_{\beta_0, \dots, \beta_p} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2$$

↑                            ↓  
Dependent Variable      Constant/Intercept  
                                ↑                            ↓  
                                Slope/Coefficient      Independent Variable

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_i X_i$$

**Y : Dependent variable**  
 **$\beta_0$  : Intercept**  
 **$\beta_i$  : Slope for  $X_i$**   
**X = Independent variable**



# Different variable selection methods

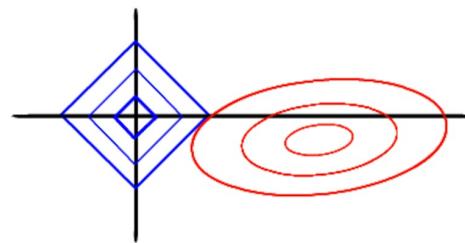
Principal Component Analysis

LASSO: least absolute shrinkage and selection operator

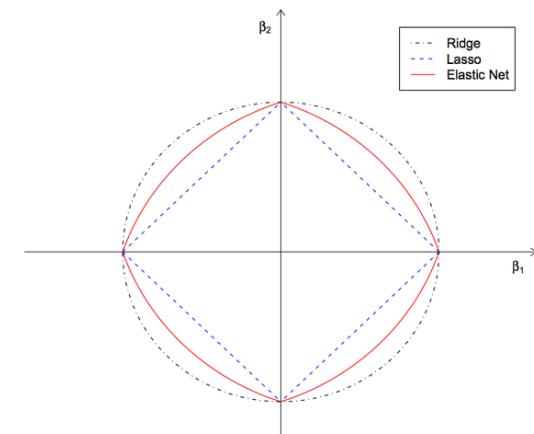
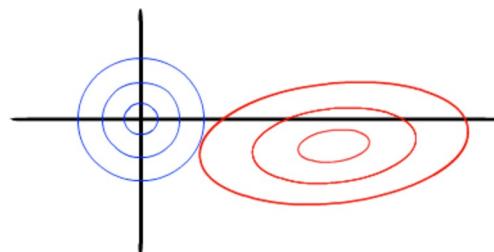
$$\arg \min_{\beta_0, \dots, \beta_p} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

The Elastic Net

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda_2 \|\mathbf{w}\|_2^2 + \lambda_1 \|\mathbf{w}\|_1$$



Ridge regression



# Applications of data aggregation in relationships

A single regression



Ryzhov, Han, and Bradic (MNSC 2016)

Multiple regressions



Li, Netessine, and Koulayev (MNSC 2018)  
Gu, and Kannan (JMR 2025)

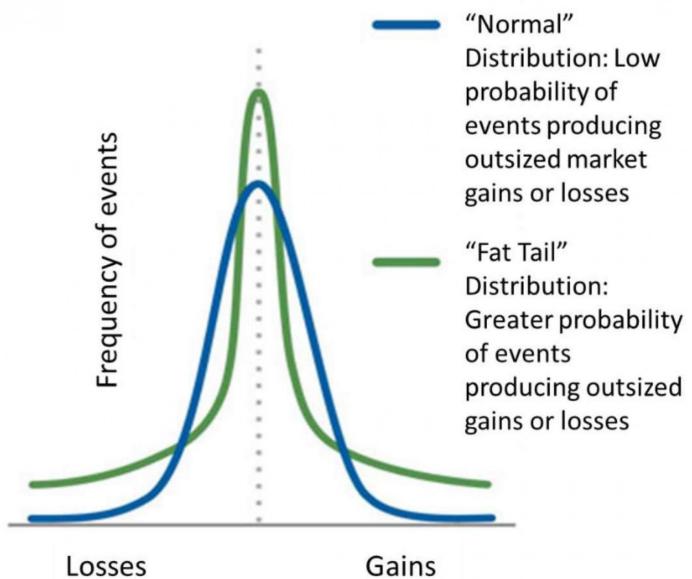


THE HONG KONG  
POLYTECHNIC UNIVERSITY  
香港理工大學

Background	One Entity	Multi Entities	Relationships	Implications
<ul style="list-style-type: none"> <li>• Prevalence of data or even big data</li> <li>• Urgency to make fast and good decisions</li> <li>• Limited cognitive resource and attention</li> </ul>	<ul style="list-style-type: none"> <li>• Basic summary statistics</li> <li>• Bias in observability of estimates and forecasts</li> <li>• Bias in observed estimates and forecasts (e.g., fat tails)</li> </ul>	<ul style="list-style-type: none"> <li>• Dimension reduction</li> <li>• Market structure analysis with structured data</li> <li>• Market structure analysis with unstructured data</li> </ul>	<ul style="list-style-type: none"> <li>• Variable selection in one linear regression</li> <li>• Variable selection in multiple linear regressions</li> </ul>	<ul style="list-style-type: none"> <li>• Optimal aggregation rules</li> <li>• Impact of reviews on product sales</li> <li>• Impact of AIGS on review activities and content consumption</li> </ul>



# Optimal aggregation when estimates are fat tailed



A simple heuristic that performs quite well across all possible scenarios.

$$\text{AMA } x = \frac{1}{2} (\text{Average } x + \text{Median } x).$$

Judges	(a) GN <sub>1</sub> samples (fat tails)					
	Bayes GN <sub>1</sub>	Bayes GN <sub>1.5</sub>	Average	Trimmed	Median	AMA
3	0%	3%	6%	—	4%	0%
5	0%	4%	12%	1%	5%	2%
10	0%	6%	20%	8%	1%	4%
20	0%	7%	26%	11%	2%	6%

Judges	(b) GN <sub>1.5</sub> samples (intermediate tails)					
	Bayes GN <sub>1</sub>	Bayes GN <sub>1.5</sub>	Average	Trimmed	Median	AMA
3	0%	0%	1%	—	11%	2%
5	1%	0%	2%	3%	12%	2%
10	3%	0%	3%	1%	8%	1%
20	5%	0%	4%	1%	11%	1%

Judges	(c) GN <sub>2</sub> samples (normal tails)					
	Bayes GN <sub>1</sub>	Bayes GN <sub>1.5</sub>	Average	Trimmed	Median	AMA
3	2%	0%	0%	—	16%	4%
5	5%	1%	0%	6%	19%	5%
10	10%	2%	0%	3%	18%	5%
20	14%	3%	0%	3%	21%	6%

Lobo, Miguel, and Dai Yao (2025). Fat tails in human judgment



# Impact of ratings on product sales and subsequent ratings

TABLE 4  
The Effects of WOM Communication on Weekly Box Office Revenue

Week (t)	1	2	3	4	5	6	7	8
CONST	7.273 (.000)**	9.346 (.000)**	6.449 (.000)**	8.171 (.000)**	6.469 (.000)**	4.981 (.015)**	.594 (.000)**	4.067 (.317)
LNSCRN <sub>t</sub>	.566 (.000)**	.401 (.006)**	.673 (.000)**	.985 (.000)**	.973 (.000)**	1.050 (.000)**	1.379 (.000)**	1.185 (.000)**
LNMSG <sub>t-1</sub>	.592 (.000)**	.345 (.000)**	.366 (.000)**	.275 (.005)**	.387 (.000)**	.169 (.0103)	-.088 (.449)	.144 (.362)
POSIPER <sub>t-1</sub>	.784 (.257)	-.300 (.674)	.812 (.133)	.186 (.603)	.342 (.312)	.045 (.910)	-.185 (.497)	.118 (.715)
NEGPER <sub>t-1</sub>	.059 (.952)	-1.369 (.092)*	-.610 (.264)	.445 (.309)	.385 (.136)	-.173 (.620)	-.265 (.309)	-.605 (.655)
LNCRITIC	.759 (.018)**	1.106 (.001)**	.894 (.020)**	.343 (.323)	.421 (.137)	.657 (.192)	.358 (.228)	.561 (.369)
CRPRO	.695 (.019)**	.398 (.235)	.688 (.076)*	.728 (.047)**	.736 (.009)**	.340 (.336)	.386 (.288)	.307 (.454)
LNNEW	-.111 (.505)	.005 (.984)	.153 (.494)	-.251 (.127)	-.136 (.423)	-.086 (.738)	.073 (.711)	-.112 (.682)
LNAGE	.146 (.760)	-.339 (.516)	.022 (.972)	-.1378 (.045)**	-.681 (.140)	-.039 (.964)	-1.510 (.097)*	.060 (.963)
Model fit F (p value)	44.87 (.000)	16.86 (.000)	22.70 (.000)	27.50 (.000)	50.37 (.000)	23.58 (.000)	69.04 (.000)	22.45 (.000)
Adjusted R <sup>2</sup>	.902	.765	.820	.855	.925	.862	.954	.905
Joint F test of WOM measures	14.01 (.000)	9.64 (.000)	10.87 (.000)	4.01 (.017)	10.48 (.000)	1.07 (.382)	.43 (.735)	.53 (.672)

\**p* < .10.

\*\**p* < .05.

Notes: The dependent variable is LNREV<sub>t</sub>. *p* values are in parentheses.

Liu, Yong (JM 2006). Word of mouth for movies: Its dynamics and impact on box office revenue



# Impact of AIGS on review activities and content consumption

Table 3: Impact of AIGS on Online Video Consumption

	DailyViews (1)	DailyLikes (2)	DailyShares (3)	DailyReviews (4)	DailyTips (5)
$Treat_i * After_t$	.142*** (.0257)	.0758*** (.0155)	.0332*** (8.51e-03)	.0294** (.0106)	.0246* (.0104)
$Followers_{it}$	1.73e-05*** (4.30e-06)	-1.25e-06 (2.59e-06)	-1.23e-05*** (1.42e-06)	-1.29e-05*** (1.77e-06)	-8.84e-06*** (1.73e-06)
$Age_{it}$	-0.0554*** (4.34e-03)	-0.0454*** (2.62e-03)	-0.0152*** (1.42e-03)	-0.0260*** (1.76e-03)	-0.0213*** (1.72e-03)
$Age_{it}^2$	-1.26e-04** (4.41e-05)	1.52e-04*** (2.66e-05)	9.63e-05*** (1.46e-05)	2.07e-04*** (1.82e-05)	1.57e-04*** (1.77e-05)
Video fixed effect	Yes	Yes	Yes	Yes	Yes
Time fixed effect	Yes	Yes	Yes	Yes	Yes
N	29,363	29,363	29,363	29,363	29,363
R <sup>2</sup>	7.4%	7.0%	3.7%	5.1%	3.6%

16

Notes: (1) Robust standard errors clustered at the video level in parentheses;

(2) \*  $p < 0.05$ , \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$



Qin, Rui, Yue Katherine Feng, and Dai Yao (2024). The impact of AI-generated summaries on video consumption

Background	One Entity	Multi Entities	Relationships	Implications
<ul style="list-style-type: none"> <li>• Prevalence of data or even big data</li> <li>• Urgency to make fast and good decisions</li> <li>• Limited cognitive resource and attention</li> </ul>	<ul style="list-style-type: none"> <li>• Basic summary statistics</li> <li>• Bias in observability of estimates and forecasts</li> <li>• Bias in observed estimates and forecasts (e.g., fat tails)</li> </ul>	<ul style="list-style-type: none"> <li>• Dimension reduction</li> <li>• Market structure analysis with structured data</li> <li>• Market structure analysis with unstructured data</li> </ul>	<ul style="list-style-type: none"> <li>• Variable selection in one linear regression</li> <li>• Variable selection in multiple linear regressions</li> </ul>	<ul style="list-style-type: none"> <li>• Optimal aggregation rules</li> <li>• Impact of reviews on product sales</li> <li>• Impact of AIGS on review activities and content consumption</li> </ul>

