

# Feature Analysis and Visualization

MIE223  
Winter 2025

## 1 Feature Analysis and Visualization

### 1.1 Exploratory Data Analysis (EDA)

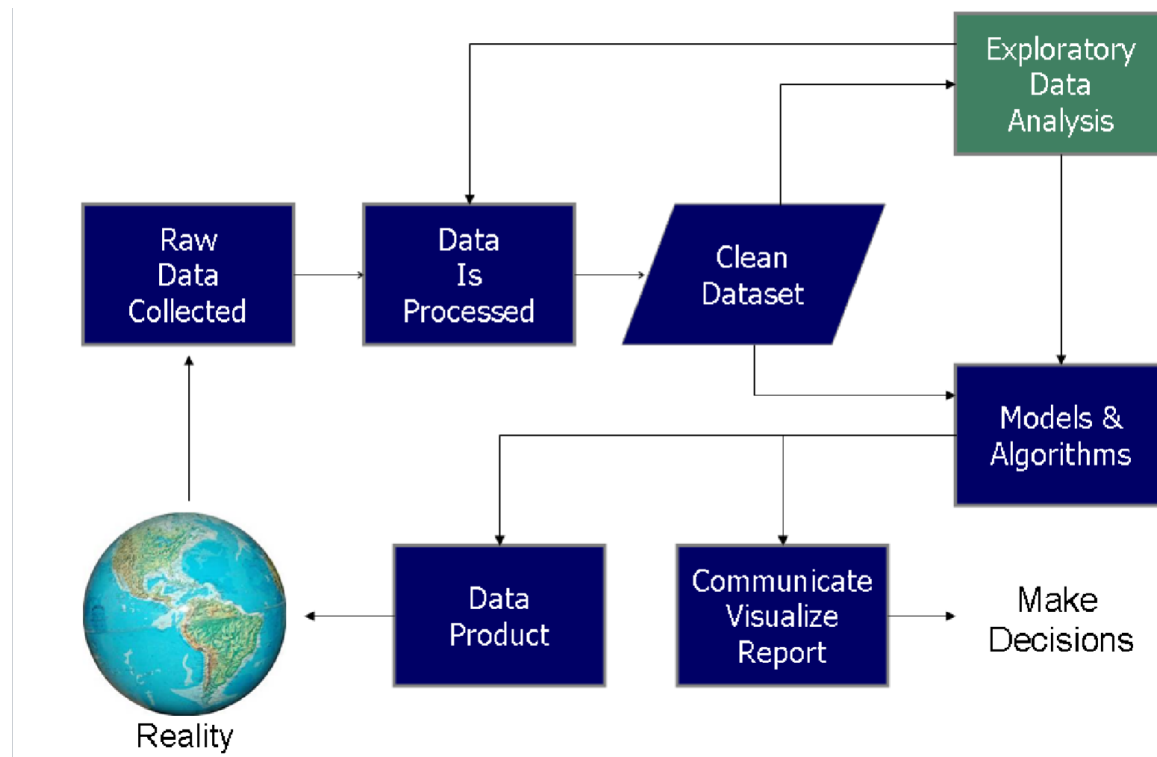
John Tukey (Bell Labs, 1970's) (Inventor of the box-plot)

- The first time corporations had a lot of digital data
- Data on semiconductor processes, networks

Observed that statistics was preoccupied with hypothesis testing

- But there were no established methodologies for **generating (data-driven) hypotheses**
- Espoused a visual methodology and a new language S (R became the free version)

### 1.2 Data Science Process



3

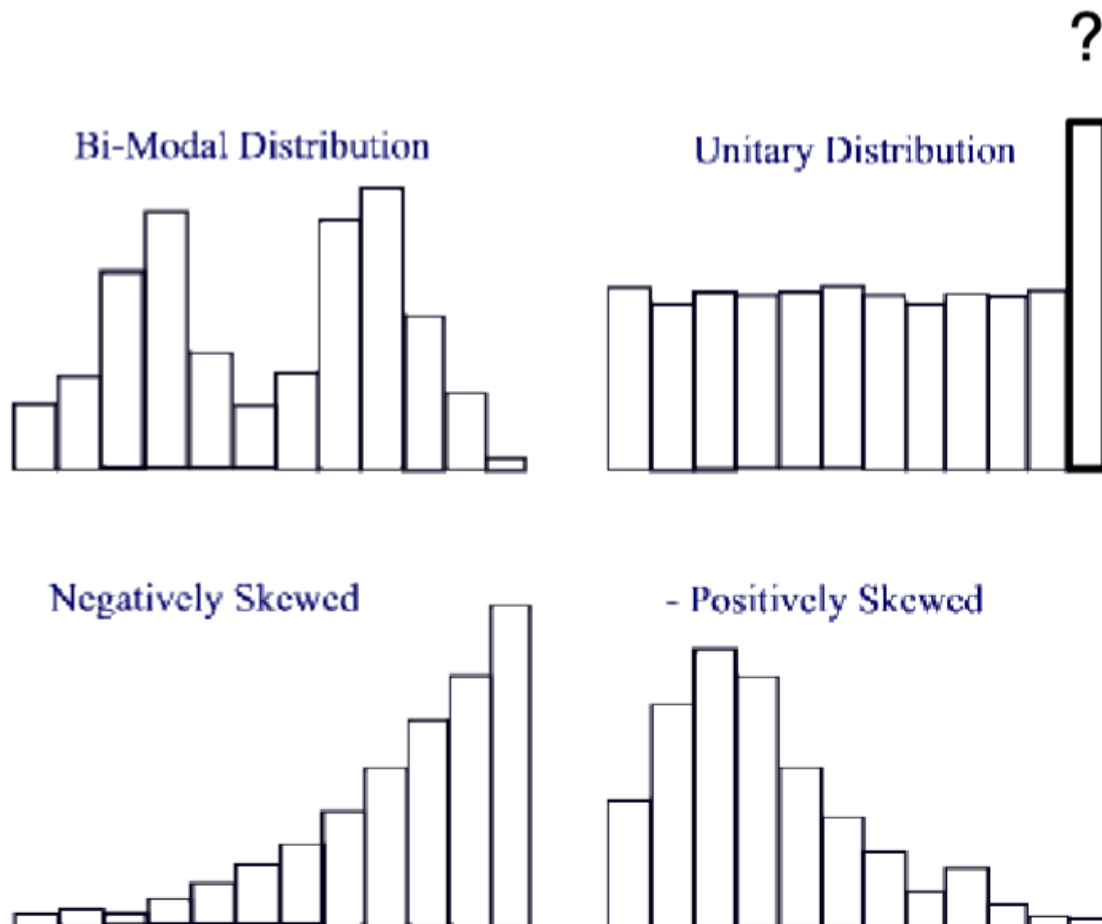
### 1.3 Data Cleaning: Brief Recap

80% of your Data Science time

Some pointers:

- Missing values – do not replace with 0 or -999

- Treat as missing
- Look for frequencies of values or outliers (-999)
  - **Histograms** invaluable
  - Examine **outliers**



## 2 Correlation

We'll be look at feature distributions independently (per column), but also relationships between features (columns). Correlation is one of the first tools you think of for looking at the relationship between two features (but has limitations)

### 2.1 Correlated variables

Two variables are correlated if changes in one variable, correspond to changes in the other  
Correlation in general refers to the statistical association between two variables

- This statistical association allows us to make estimates for the one variable based on the value of the other

- While we typically think of linear relationships between variables, nonlinear might exist too

## 2.2 Linear correlation

- Two variables  $x$  and  $y$  are said to be linearly correlated if their relationship can be described by the following equation:

$$y = \alpha + \beta x, \beta \neq 0 \quad (1)$$

- This relationship will not be exact
- There will be an error associated with it, and hence, in reality:
- $\epsilon$  is the error term for the relationship

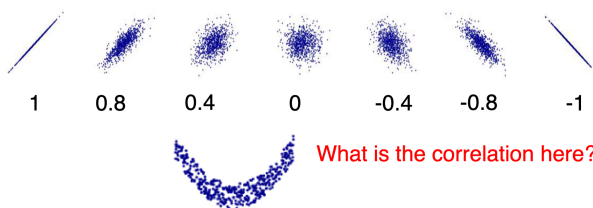
$$y = \alpha + \beta x + \epsilon, \beta \neq 0 \quad (2)$$

## 2.3 Pearson correlation coefficient

- In order to test the linear association between two variables  $x$  and  $y$  we can use the Pearson correlation coefficient  $r_{xy}$

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

- Pearson correlation is in the range  $[-1, 1]$ 
  - 1: perfect/strong and positive linear correlation
  - -1: perfect/strong and negative linear correlation
  - 0: no linear correlation
- It is very crucial to understand that this correlation coefficient can only examine linear associations between two variables



**Note 1.** • Pearson correlation coefficient is sensitive to outliers

- It is not robust to outliers
- It is not a good measure of correlation for non-linear relationships

### 3 Feature Analysis: Probability Recap

This is mathematical so we first we need to make sure everyone has the same understanding and intuition for notation

#### 3.1 Random Variables

- Random variable (RV) denoted by uppercase letter (e.g.  $X$ )
- RVs take value assignments  $X = x$  where  $X$  is a
  - Discrete RV if  $x$  is in a countable set (binary:  $x \in \{0,1\}$ ; or dice)
  - Continuous RV if  $x$  is in an uncountable set (real:  $x \in \text{Reals}$ )
- Write  $x \in X$  for possible value assignments of  $X$
- For all  $x$ ,  $P(X = x) \in [0,1]$
- $P(X = x)$  is a proper distribution
  - Discrete RV:
$$\sum_{x \in X} P(X = x) = 1$$
  - Continuous RV:
$$\int_{x \in X} P(X = x) dx = 1$$
- Write  $P(x)$  for  $P(X=x)$ , write  $P(X)$  for full distribution
- Representing probability distributions
  - Discrete (finite) RV: tabular
  - Continuous (infinite) RV: function

#### 3.2 Joint Distributions on RVs

Aliens in your backyard

- Aliens in your backyard

$P(R,C) =$

$R$	$C$	$Pr$
1	1	0
1	2	.4
2	1	.2
2	2	.4

	$C=1$	$C=2$
$R=1$	0	.4
$R=2$	.2	.4

- $P(R=2,C=2) = .4$
- $P(R=2) = \sum_{c \in \{1,2\}} P(R=2,C=c) = .6$
- $P(R=1|C=2) = P(R=1,C=2)/P(C=2) = .5$

Marginalize over C

Condition on C=2

Example from Andrew  
Moore @ CMU/Google

### 3.3 Rules of Probability

- Joint and conditional distributions:

$$P(X,Y) = P(X|Y)P(Y) = P(Y|X)P(X) \quad (3)$$

- Marginalization:

$$P(X) = \sum_y P(X,Y) \quad (4)$$

- Conditional probability & Bayes rule:

$$P(X|Y) = \frac{P(Y|X)P(X)}{P(Y)} \quad (5)$$

### 3.4 Manipulating Distributions

- Sometimes we don't just want  $P(R=1,C=2) = .4$
- We want to work with full distributions  $P(R,C)$
- How to apply previous rules to full distributions?
- easy, just do once for each case and store in table
- Marginalization:

$$\sum_b P(A, b) = P(A)$$

$\sum_b$ 

$A$	$B$	$Pr$
0	0	.1
0	1	.3
1	0	.2
1	1	.4

 $=$ 

$A$	$Pr$
0	.4
1	.6

•

- Binary Multiplication

$$P(A) \cdot P(B|A) = P(A, B)$$

<table><tr><th><math>A</math></th><th><math>Pr</math></th></tr><tr><td><math>0</math></td><td><math>.7</math></td></tr><tr><td><math>1</math></td><td><math>.3</math></td></tr></table>	$A$	$Pr$	$0$	$.7$	$1$	$.3$	$\cdot$	<table><tr><th><math>A</math></th><th><math>B</math></th><th><math>Pr</math></th></tr><tr><td><math>0</math></td><td><math>0</math></td><td><math>.1</math></td></tr><tr><td><math>0</math></td><td><math>1</math></td><td><math>.9</math></td></tr><tr><td><math>1</math></td><td><math>0</math></td><td><math>.2</math></td></tr><tr><td><math>1</math></td><td><math>1</math></td><td><math>.8</math></td></tr></table>	$A$	$B$	$Pr$	$0$	$0$	$.1$	$0$	$1$	$.9$	$1$	$0$	$.2$	$1$	$1$	$.8$	$=$	<table><tr><th><math>A</math></th><th><math>B</math></th><th><math>Pr</math></th></tr><tr><td><math>0</math></td><td><math>0</math></td><td><math>.07</math></td></tr><tr><td><math>0</math></td><td><math>1</math></td><td><math>.63</math></td></tr><tr><td><math>1</math></td><td><math>0</math></td><td><math>.06</math></td></tr><tr><td><math>1</math></td><td><math>1</math></td><td><math>.24</math></td></tr></table>	$A$	$B$	$Pr$	$0$	$0$	$.07$	$0$	$1$	$.63$	$1$	$0$	$.06$	$1$	$1$	$.24$
$A$	$Pr$																																							
$0$	$.7$																																							
$1$	$.3$																																							
$A$	$B$	$Pr$																																						
$0$	$0$	$.1$																																						
$0$	$1$	$.9$																																						
$1$	$0$	$.2$																																						
$1$	$1$	$.8$																																						
$A$	$B$	$Pr$																																						
$0$	$0$	$.07$																																						
$0$	$1$	$.63$																																						
$1$	$0$	$.06$																																						
$1$	$1$	$.24$																																						

•

## 4 Feature Analysis

Frequency, Correlation, (Pointwise) Mutual Information, Conditional Entropy

### 4.1 Feature Analysis: Frequency

- Always produce summary frequency statistics
- Always look at samples of data

Analysis of Twitter Data from 2013/14  
crawled according to topical hashtags

Have your  
own software  
export these  
tables...  
why?

#Unique Features				
From	Hashtag	Mention	Location	Term
95,547,198	11,183,410	411,341,569	58,601	20,234,728

	#Unique Features									
	From	Hashtag	Mention	Location	Term	From	Hashtag	Mention	Location	Term
#TrainHashtags	58	98	126	12	49	28	31	51	52	29
#TestHashtags	36	63	81	5	29	16	19	19	33	17
#TopicalTerms	51,053	2,927,719	500,380	8,362	408,304	163,800	3,000,058	2,800,058	2,013,117	382,537
Sample Hashtags	#toscanchampion	#asteroids	#workday	#randall	#jordanerick	#robwilliams	#policebrutality	#earthquake	#ebola	#veividev
	#newadglovis	#astronauts	#wessexer	#unfathom	#chiderofpita	#jennamela	#michelebrown	#stern	#vms	#gettyvide
	#wastickon	#astefine	#tita	#small	#super	#rijanarrows	#picksell	#summi	#sacato	#astefine
	#concomeniti	#spacecraft	#realmaid	#mham	#combthre	#randella	#millerwood	#ablood	#chickenpot	#omo
	#mimnews	#elkscope	#beckham	#nuclearpower	#vms	#paulwalker	#wigganlan	#murrenkatrim	#teplague	#gemmang

•

### 4.2 Feature (In)dependence

Often we want to know whether one feature column predicts another (bivariate feature analysis)

- What do you see here?
- How would quantify / visualize / automatically detect this?

Person ID	Gender	Age	Party
1	M	25	Dem
2	M	35	Rep
3	M	45	Rep
4	F	25	Dem
5	F	35	Dem
6	F	45	Rep
7	F	55	Rep

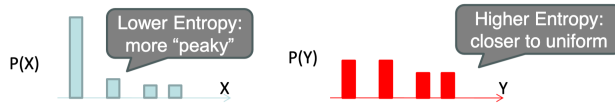
If numeric, could just look at correlation or  $R^2$ , but here we have discrete labels

### 4.3 Feature Independence

- X is independent of Y means that knowing Y does not change our belief about X
- $p(X|Y=y) = p(X)$
- In words: knowing  $Y=y$  has no impact on our belief in the distribution of X
- Occurs if:  $p(X=x, Y=y) = p(X=x) \cdot p(Y=y)$

### 4.4 Information Theory

- $P(X)$  encodes our uncertainty about X
- Some variables are more uncertain than others



- How can we quantify this intuition?
- Entropy: average number of bits required to encode X

$$H_p(X) = E\left[\log \frac{1}{p(x)}\right] = \sum_x P(x) \log \frac{1}{P(x)} = -\sum_x P(x) \log P(x)$$

Higher entropy is more uncertain, lower entropy is more certain!

- In short: low bits, low entropy  $\rightarrow$  low uncertainty
- We can define conditional entropy (CE) similarly

$$H_p(X|Y) = E\left[\log \frac{1}{p(x|y)}\right] = H_p(X, Y) - H_p(Y)$$

- i.e. once Y is known, we only need  $H(X, Y) - H(Y)$  bits
- Higher CE is more uncertain (y is less predictive), lower CE is more certain (y is more predictive)

## 4.5 (Pointwise) Mutual Information

**Note 2.** Mutual information is just a test of independence. It only applies to discrete variables. Correlation does not apply to binary variables. If X and Y are non-linear, MI will capture that unlike correlation coefficient.

- Mutual Information

$$I(X; Y) = \sum_{y \in Y} \sum_{x \in X} p(x, y) \log \left( \frac{p(x, y)}{p(x) p(y)} \right)$$

- For discrete random variables: can also be generalized to continuous case
- PMI: target specific values x and y (pointwise)

$$\text{pmi}(x; y) \equiv \log \frac{p(x, y)}{p(x)p(y)} = \log \frac{p(x|y)}{p(x)} = \log \frac{p(y|x)}{p(y)}$$

- Lower MI and PMI → More independent, Higher MI and PMI → More dependent

## 4.6 Discrete Feature Analysis Summary

If we look at two feature columns X and Y, can analyze:

- MI(X,Y): test independence (fix neither variable)
  - MI=0 → independent,
  - The higher the value, the more dependence
  - Chi2(X,Y) / Chi-squared(X,Y): independence test similar to MI
- CE(X—Y=b): test how well Y=b predicts X (fix Y)
  - Low conditional entropy when X is well-predicted, why?
  - CE=0 → perfect prediction! CE=1 → complete noise!
- PMI(X=a,Y=b): test if X=a, Y=b dependent (fix X and Y)
  - Hard to interpret absolute numerical meaning (0 is still independent)
  - Primarily useful when ranking features... why?

## 4.7 Predictive Feature Analysis: Mutual Information (MI) or Pointwise MI

E.g., top 5 term features for Twitter topics – Use mutual information to know what is predictive / important



## 4.8 Mean, Median, and Power Laws

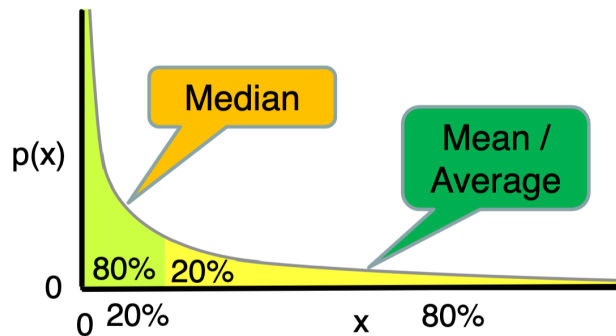
- A lot of human data is power law (income, #friends)
- Discrete power law (Zipf) for freq  $f$ :

- $p(f) = \alpha f^{-1-1/s}$

- Continuous power law:

- $P(X > x) \sim L(x)x^{-\alpha+1}$

- E.g., Probability mass of population with different income levels, or #friends on Facebook

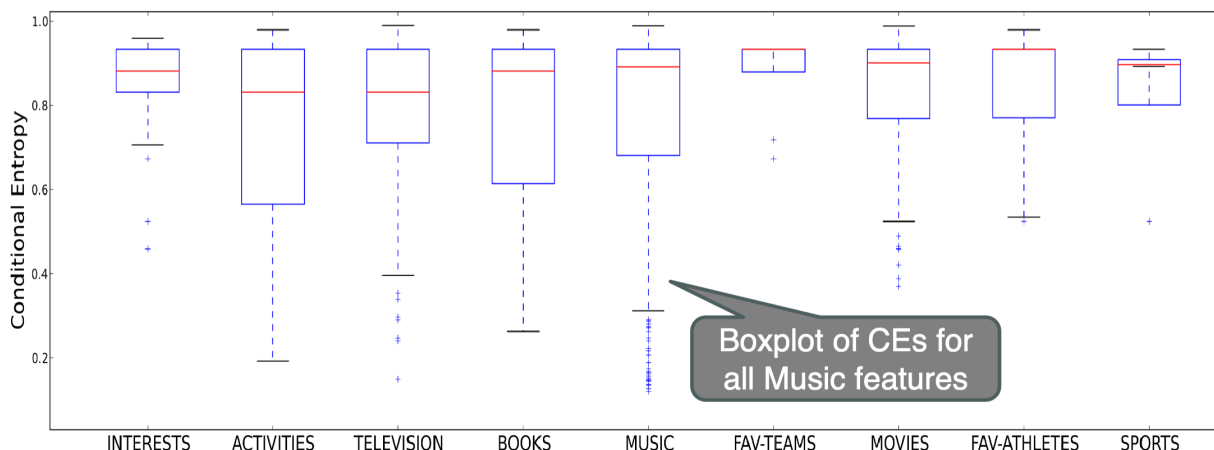


- Not symmetric like Gaussian, mean  $\neq$  median
- I.e., the average case is much closer to the max case!
- 80/20 rule: top 20% of values account for 80% of mass

## 4.9 Box Plots

Most data is not Gaussian, often power law or log-Normal

How well can I predict a user's liking behavior conditioned on observing a feature, e.g., you like "U2" or "Breaking Bad"



An example where most informative = mean informative (but not median)... why?

## Median Informative

Median Informative Favourites by Category		
Movies	Music	Television
Forrest Gump	John Lennon	Futurama
Pretty Woman	U2	Star Trek
Napoleon Dynamite	AC/DC	The Trap Door
Harry Potter	The Smashing Pumpkins	Drawn Together
Toy Story 3	Gotye	Sherlock(Official)
The Godfather	The Rolling Stones	Hitchhiker's Guide to the Galaxy
Mulan	All Access	Buffy The Vampire Slayer
How to Train Your Dragon	Steve Aoki	South Park

## Most Informative

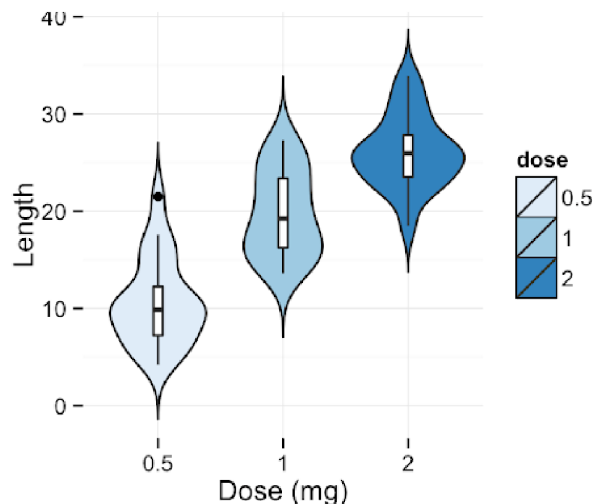
Most Informative Favourites by Category		
Movies	Music	Television
Billy Madison	Avascular Necrosis	Metalocalypse
Team America: World Police	Tortured	Beast Wars
Pan's Labyrinth	Elysian	Hey Arnold!
Pirates of the Caribbean	Anno Domini	Sherlock
Aladdin	Darker Half	Hey Hey It's Saturday
Starship Troopers	Hellbringer	Neil Buchanan and Art Attack!
Happy Gilmore	Johnny Roadkill	Breaking Bad
Timon and Pumbaa	Aeon of Horus	Red vs. Blue

- Median favorites were largely generic
- Most informative (max  $\approx$  avg) were largely specialized

### 4.10 Violin Plots

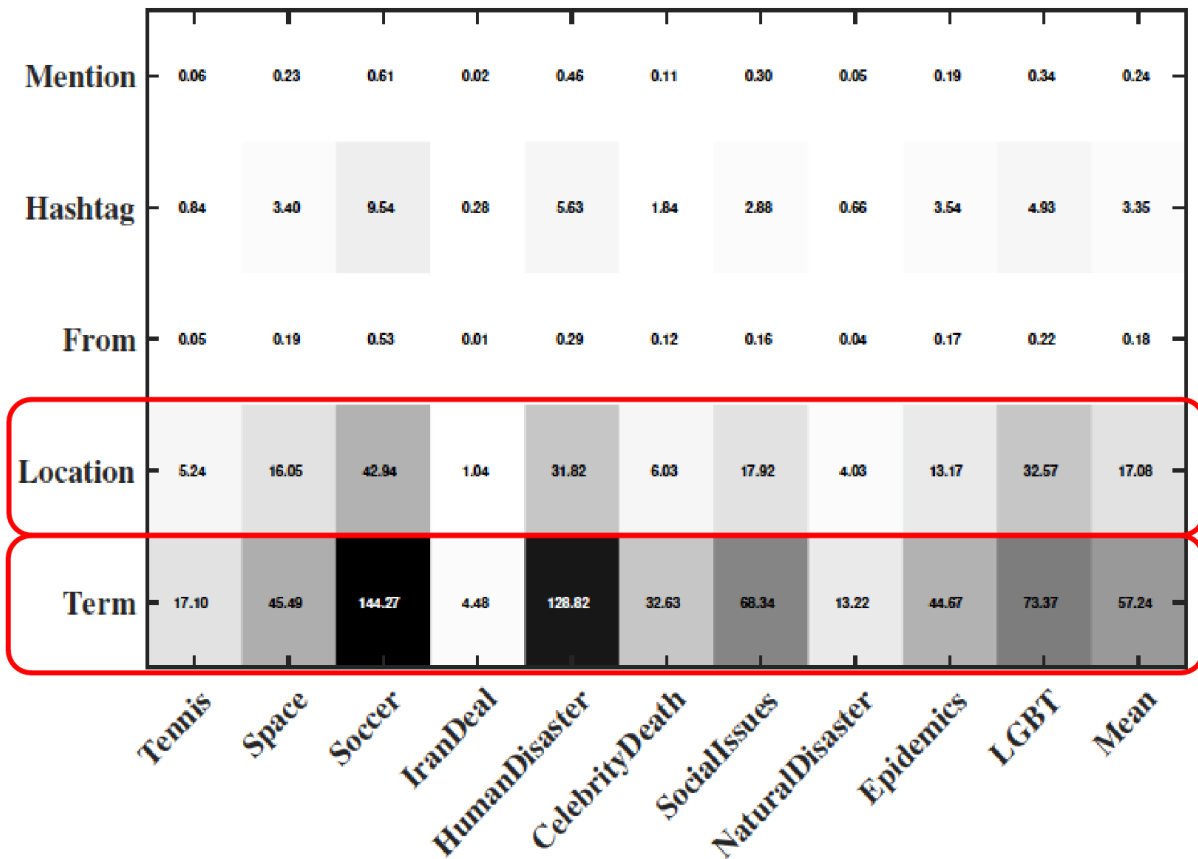
Compactly show full distributions side-by-side. Like compact histogram, can show bimodality unlike Box-plot

Plot of tooth length vs. vitamin C dose



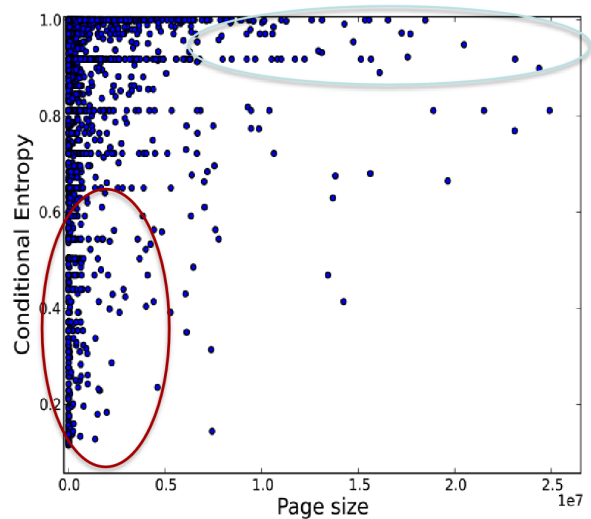
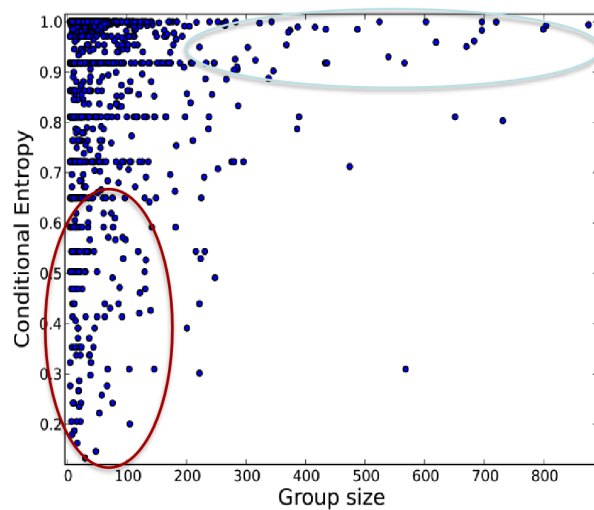
### 4.11 Discrete Heatmaps

Better than a table... visualize numbers!

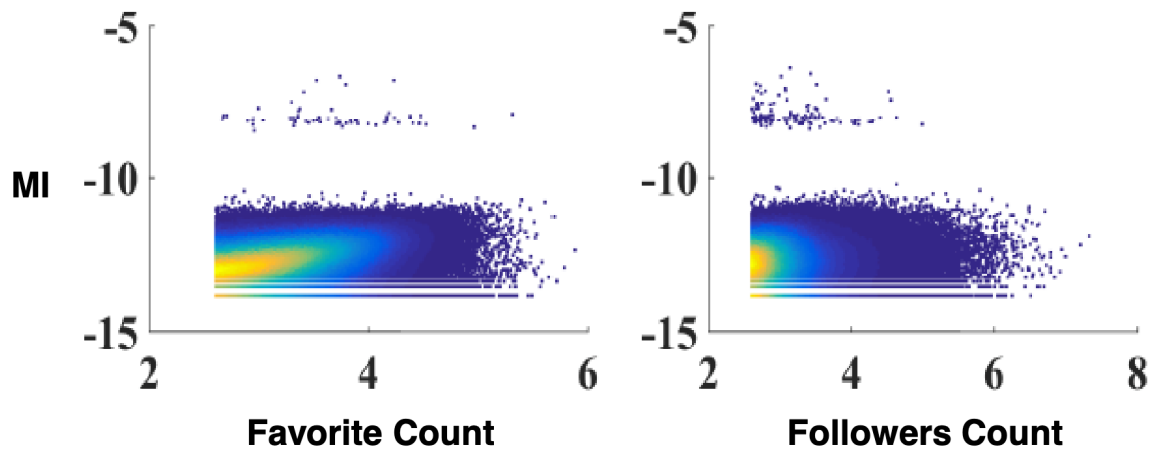


#### 4.12 Scatterplots: don't need linear relationship to be informative!

- Large groups tend not to be predictive
- Most informative groups were small



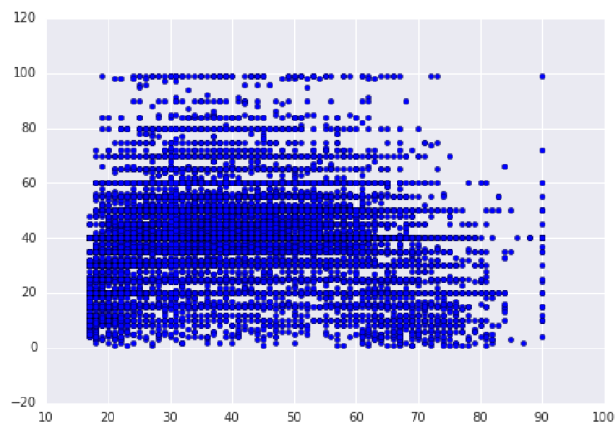
### 4.13 Heatmaps / Density Plots



Important to use when points in a scatterplot are dense and do not reflect density

### 4.14 Bubble Plots

Sometimes scatterplots not dense enough for a heatmap



- And may have a third data dimension
- So don't want density=color



Enter the bubble plot

- X-axis: age
- Y-axis: hours worked per week
- Bubble size is relative mass (density in a heatmap)
- Color is income (purple=higher)