



### Learning Objectives

- 1. Identify *properties* of data and images
- 2. Decide how to encode data using visual attributes/channels
- 3. Define concepts of *expressiveness* and *effectiveness*
- 4. Develop *automated chart design* algorithm



TODAY











# **DATA MODELS & CONCEPTUAL MODELS**

### Data models are formal descriptions

**Math:** Sets with operations on them **Examples:** integers with +, - and × operators reals/floats with +, -, × and ÷

### **Conceptual models** are mental constructions

Include semantics and support reasoning

### Examples (data vs. conceptual)

1D floats vs. temperature 3D tuple of floats vs. spatial location in 3D

	competitorname	chocolate	fruity	caramel	peanutyalmondy	nougat	crispedricewafer	hard	bar	pluribus	sugarpercent	pricepercent	winpercent
1	100 Grand	1	0	1	0	0	1	0	1	0	.73199999	.86000001	66.971725
•	3 Musketeers	1	0	0	0	1	0	0	1	0	.60399997	.51099998	67.602936
1	One dime	0	0	0	0	0	0	0	0	0	.011	.116	32.261086
5	One quarter	0	0	0	0	0	0	0	0	0	.011	.51099998	46.116505
6	Air Heads	0	1	0	0	0	0	0	0	0	.90600002	.51099998	52.341465
7	Almond Joy	1	0	0	1	0	0	0	1	0	.465	.76700002	50.347546
8	Baby Ruth	1	0	1	1	1	0	0	1	0	.60399997	.76700002	56.914547
9	Boston Baked Beans	0	0	0	1	0	0	0	0	1	.31299999	.51099998	23.417824
0	Candy Corn	0	0	0	0	0	0	0	0	1	.90600002	.32499999	38.010963
1	Caramel Apple Pops	0	1	1	0	0	0	0	0	0	.60399997	.32499999	34.517681
2	Charleston Chew	1	0	0	0	1	0	0	1	0	.60399997	.51099998	38.975037
3	Chewey Lemonhead Fruit Mix	0	1	0	0	0	0	0	0	1	.73199999	.51099998	36.017628
I	string	bool	bool	bool	haal	bool	haal	bool	bo	ol bool	float	float	float
- 12	string			1 2000.			1000	2000.			noat	Inoat	Inoat

### **CONCEPTUAL MODEL**

Header	Description
chocolate	Does it contain chocolate?
fruity	Is it fruit flavored?
caramel	Is there caramel in the candy?
peanutalmondy	Does it contain peanuts or almonds?
nougat	Does it contain nougat?
crispedricewafer	Does it contain crisped rice or cookies?
hard	Is it a hard candy?
bar	Is it a candy bar?
pluribus	Is it one of many candies in a bad?
sugarpercent	The percentile of sugar (across dataset)
pricepercent	The unit price percentile (across dataset)
winpercent	The overall win percentage in 269K contests

https://github.com/fivethirtyeight/data/tree/master/candy-power-ranking

17

# **CONCEPTUAL MODEL**

Header	Description
chocolate	Does it contain chocolate?
fruity	Is it fruit flavored?
caramel	Is there caramel in the candy?
peanutalmondy	Does it contain peanuts or almonds?
nougat	Does it contain nougat?
crispedricewafer	Does it contain crisped rice or cookies?
hard	Is it a hard candy?
bar	Is it a candy bar?
pluribus	Is it one of many candies in a bad?
sugarpercent	The percentile of sugar (across dataset)
pricepercent	The unit price percentile (across dataset)
winpercent	The overall win percentage in 269K contests

Domain specific understanding of the data

Supports analysis and reasoning



### **DATA TYPES**

**N - Nominal (labels)** Fruits: Apples, oranges, ...

Operations: =, ≠

O - Ordered Quality of eggs: Grade AA, A, B Operations: =, ≠, <, >

#### **Q** - Interval (location of zero arbitrary)

Dates: Jan, 19, 2016; Loc.: (LAT 33.98, LON -118.45) Like a geometric point. Cannot compare directly Only differences (i.e. intervals) may be compared Operations =,  $\neq$ , <, >, -

#### **Q** - Ratio (location of zero fixed)

Physical measurement: Length, Mass, ... Counts and amounts Like a geometric vector, origin is meaningful Operations: =,  $\neq$ , <, >, -,  $\div$ 

# NOMINAL, ORDINAL, QUANTITATIVE

Header	Description	
competitorname	Name of candy	Ν
chocolate	Does it con <mark>tain chocolate?</mark>	N (maybe O)
fruity	ls it fruit fla <mark>vored?</mark>	N (maybe O)
caramel	Is there cara <mark>mel in the candy?</mark>	N (maybe O)
peanutalmondy	Does it contain peanuts or almonds?	N (maybe O)
nougat	Does it con <mark>tain nougat?</mark>	N (maybe O)
crispedricewafer	Does it contain crisped rice or cookies?	N (maybe O)
hard	ls it a hard <mark>candy?</mark>	N (maybe O)
bar	ls it a candy bar?	N (maybe O)
pluribus	Is it one of many candies in a bad?	N (maybe O)
sugarpercent	The percentile of sugar (across dataset)	Q-Ratio
pricepercent	The unit price percentile (across dataset)	Q-Ratio
winpercent	The overall win percentage in 269K contests	Q-Ratio

# DATA TYPES

### DIMENSIONS

Dimensions are often the **independent** variables

*Dimensions* contain **qualitative values that describe the data item** (such as names, dates, or geographical data)

### **MEASURES**

Measures are often the **dependent** variables

*Measures* contain numeric, **quantitative** values that you can measure *in the experiment*. Measures can be aggregated (sum, count, average, std. deviation).

1	competitorname	chocolate	fruity	caramel	peanutyalmondy	nougat	crispedricewafer	hard	bar	pluribus	sugarpercent	pricepercent	winpercent
2	100 Grand	1	0	1	0	0	1	0	1	0	.73199999	.86000001	66.971725
3	3 Musketeers	1	0	0	0	1	0	0	1	0	.60399997	.51099998	67.602936
4	One dime	0	0	0	0	0	0	0	0	0	.011	.116	32.261086

NOTE: Distinction is not strict. The same variable may be treated either way depending on the task

21

# **DIMENSION OR MEASURE**

Header	Description
competitorname	Name of candy
chocolate	Does it contain chocolate?
fruity	Is it fruit flavored?
caramel	Is there caramel in the candy?
peanutalmondy	Does it contain peanuts or almonds?
nougat	Does it contain nougat?
crispedricewafer	Does it contain crisped rice or cookies?
hard	Is it a hard candy?
bar	ls it a candy bar?
pluribus	Is it one of many candies in a bad?
sugarpercent	The percentile of sugar (across dataset)
pricepercent	The unit price percentile (across dataset)
winpercent	The overall win percentage in 269K contests

### **DIMENSION OR MEASURE**

Header	Description
competitorname	Name of candy
chocolate	Does it contain chocolate?
fruity	Is it fruit flavored?
caramel	Is there caramel in the candy?
peanutalmondy	Does it contain peanuts or almonds?
nougat	Does it contain nougat?
crispedricewafer	Does it contain crisped rice or cookies?
hard	Is it a hard candy?
bar	Is it a candy bar?
pluribus	Is it one of many candies in a bad?
sugarpercent	The percentile of sugar (across dataset)
pricepercent	The unit price percentile (across dataset)
winpercent	The overall win percentage in 269K contests

https://github.com/fivethirtyeight/data/tree/master/candy-power-ranking

23

### **DIMENSION OR MEASURE**

Header	Description
competitorname	Name of candy
chocolate	Does it contain chocolate?
fruity	Is it fruit flavored?
caramel	Is there caramel in the candy?
peanutalmondy	Does it contain peanuts or almonds?
nougat	Does it contain nougat?
crispedricewafer	Does it contain crisped rice or cookies?
hard	Is it a hard candy?
bar	Is it a candy bar?
pluribus	Is it one of many candies in a bad?
sugarpercent	The percentile of sugar (across dataset)
pricepercent	The unit price percentile (across dataset)
winpercent	The overall win percentage in 269K contest

			А	В	С	D	E
	NCUC DATA	1	year	age	marst	sex	people
U.J. LE	NJUJ DATA	2	1850	0	0	1	1483789
		3	1850	0	0	2	1450376
		4	1850	5	0	1	1411067
People Count:	# of people in subgroup	5	1850	5	0	2	1359668
Year:	1850 – 2000 (every decade)	6	1850	10	0	1	1260099
<b>A a a</b> :	0 00	-	1850	10	0	2	1077122
Aye.	0 - 30+	9	1850	15	0	2	1110619
Sex:	Male, Female	10	1850	20	0		1017281
Marital Status	Single Married Divorced	11	1850	20	0	2	1003841
maritar Status.	Single, Married, Divorced,	12	1850	25	0	1	862547
		13	1850	25	0	2	799482
2348 data point	ts.	14	1850	30	0	1	730638
		15	1850	30	0	2	639636
		16	1850	35	0	1	588487
		17	1850	35	0	2	505012
		18	1850	40	0	1	475911
		19	1850	40	0	2	428185
		20	1850	45	0	1	384211
		21	1850	45	0	2	341254
		22	1850	50	0	1	321343
		23	1850	50	0	2	286580
		24	1850	55	0	1	194080
		25	1850	55	0	2	187208
		115	- ioch	en.			1 / 4075

			Α	В	С	D	E
CENICII		1	year	age	marst	sex	people
LEINDU	J IN, U, N	2	1850	0	0	1	1483789
		3	1850	0	0	2	1450376
Decula Count	O-Ratio	4	1850	5	0	1	1411067
People Count:	Q-Itatio	5	1850	5	0	2	1359668
Year:	Q-Interval	0	1850	10	0	1	1260099
Ado.	O-Ratio	7	1850	10	0	2	1077139
nge.	Q Hutto	9	1850	15	0	2	1110619
Sex:	Ν	10	1850	20	0	1	1017281
Marital Status	Ν	11	1850	20	0	2	1003841
lantai Status.		12	1850	25	0	1	862542
		13	1850	25	0	2	79948
		14	1850	30	0	1	730638
		15	1850	30	0	2	63963
		16	1850	35	0	1	58848
		17	1850	35	0	2	505012
		18	1850	40	0	1	47591
		19	1850	40	0	2	42818:
		20	1850	45	0	2	341254
		22	1850	50	0	1	321343
		23	1850	50	0	2	286580
		24	1850	55	0	1	194080
		25	1850	55	0	2	187208
		26	1050	60	0	1	174076

				A	В	C	D	E
CENICII		ллелс	1	year	age	marst	sex	people
<b>LENDU</b>	$\mathbf{D}$	IVIEAJ.	2	1850	0	0	1	1483789
			3	1850	0	0	2	1450376
	Moacuro		4	1850	5	0	1	1411067
People Count:	weasure		5	1850	5	0	2	1359668
Year:	Dimension		6	1850	10	0	1	1260099
<b>A</b> <i>m</i> <b>a :</b>	Moacuro		/	1850	10	0	2	1216114
Age:	Wedsure		0	1850	15	0	1	1110610
Sex:	Measure		10	1850	20	0	2	1017281
Marital Statuc	Moasuro		11	1850	20	0	2	1003841
iviaritai Status.	weasure		12	1850	25	0	- 1	862547
			13	1850	25	0	2	799482
			14	1850	30	0	1	730638
			15	1850	30	0	2	639636
			16	1850	35	0	1	588487
			17	1850	35	0	2	505012
			18	1850	40	0	1	475911
			19	1850	40	0	2	428185
			20	1850	45	0	1	384211
			21	1850	45	0	2	341254
			22	1850	50	0	1	321343
			23	1850	50	0	2	286580
			24	1850	55	0	1	194080
			25	1920	22	0	2	19/208





# RELATIONAL ALGEBRA [Codd 1970] / SQL

Projection (SELECT) – select a set of columns

select day, stock

day	stock	price	day	stock
10/3	AMZN	957.10	10/3	AMZN
10/3	MSFT	74.26	10/3	MSFT
10/4	AMZN	965.45	10/4	AMZN
10/4	MSFT	74.69	10/4	MSFT

### RELATIONAL ALGEBRA [Codd 1970] / SQL

#### Selection (WHERE) - filter rows

select \* where price > 100

day	stock	price				
10/3	AMZN	957.10		day	stock	price
10/3	MSFT	74.26	$\rightarrow$	10/3	AMZN	957.10
10/4	AMZN	965.45		10/4	AMZN	965.45
10/4	MSFT	74.69				

31

### RELATIONAL ALGEBRA [Codd 1970] / SQL

Sorting (ORDER BY) – order records

select \* order by stock

day	stock	price		day	stock	price
10/3	AMZN	957.10		10/3	AMZN	957.10
10/3	MSFT	74.26	$\rightarrow$	10/4	AMZN	965.45
10/4	AMZN	965.45		10/3	MSFT	74.26
10/4	MSFT	74.69		10/4	MSFT	74.69

# RELATIONAL ALGEBRA [Codd 1970] / SQL

#### Aggregation (GROUP BY, SUM, MIN, ...)

select stock min(price) group by stock



33

### RELATIONAL ALGEBRA [Codd 1970] / SQL

Combination (JOIN) multiple tables together

day	sto	ck	price
10/3	AM	ZN	957.10
10/3	MS	FT	74.26
10/4	AM	ZN	965.45
10/4	MS	FT	74.69
			•
STOCK	(		min
AMZN	V		957.10

	day	stock	price	min
	10/3	AMZN	957.10	957.10
-	10/3	MSFT	74.26	74.26
	10/4	AMZN	965.45	957.10
	10/4	MSFT	74.69	74.26

select t.day, t.stock, t.price, a.min
from table as t, aggregate as a
where t.stock = a.stock



# **CLASS PARTICIPATION REQUIREMENTS**

Complete required readings and notebooks before class

Attend class and be a part of the in-class discussion

Post at least 1 discussion substantive comment/question per week

Due by 8pm the following Sunday 1 free pass for the quarter

> Class home page https://magrawala.github.io/cs448b-fa24/

# **READING/NOTEBOOK/LECTURE RESPONSES**

#### Good responses typically exhibit one or more

**Critiques** of arguments made in the papers/lectures **Analysis** of implications or future directions for ideas in readings/lectures **Insightful questions** about the readings/lectures

#### **Responses should not be summaries**

Should be substantive (1-2 paragraphs is typical)



# **ASSIGNMENT 1: VISUALIZATION DESIGN**

Due Mon 9/30 6am!

### Design a static visualization for a data set

You must choose the message you want to convey. What question(s) do you want to answer? What insight do you want to communicate?

#### Updated submission guidelines

Make public Google Slide with your visualization Put your name, suid and write-up in notes Submit public link to slide on Canvas

#### Data: Stanford Olympic Medals

Stanford Athletics recently published a dataset on the Olympic Medal History of Stanford students and alumni. We have extracted and wrangled a data table containing information Olympic medals won by Stanford students since 1912. Our data contains the following information:

#### Number of records: 335

- Data fields:
- Athlete Name: Name of athlete. Host City: Name of city that hosted the Olympics
- Host Country: Name of country that hosted the Olympics.
- Year: Year the summer Olympics took place.
- Athlete Team Country: Country whose team the athlete represented
- Sport: Name of sport athlete competed in.
- Event: Name of event athlete competed in ("-" indicates event name is same as name of sport). Medal: Type of medal won (ties are indicated in parentheses).

The extracted dataset is available in csv format: olympic\_athletes-wrangled-2024.csv





# **CODING INFORMATION IN POSITION**



- 1. A, B, C are distinguishable
- 2. Three points are colinear: B between A and C
- 3. BC is twice as long as AB
- ... Encode quantitative variables

"Resemblance, order and proportional are the three signfields in graphics." - Bertin



BERTI	N'S "LE	VE	LS	OF	ORGANIZATION"
	Position	Ν	0	Q	N Nominal O Ordered
	Size	N	0	Q	Q Quantitative
	Value	Ν	0	Q	Note: $\mathbf{Q} \subset \mathbf{O} \subset \mathbf{N}$
	Texture	N	0		
	Color	Ν			
	Orientation	Ν			
	Shape	Ν			















TECHNOLOGY SOFTWARE - INFRASTRUCTURE		SEMICONDUC	CTORS			COMMUNICATION SERVICES		CONSUMER DEF DISCOUNT STORES		HEALTHCARE DRUG MANUFACTUR	ERS - GENERAL		HEALTHC	ARE PLA
	ORCL -3.16%		<b>NV</b> -0.7	<b>DA</b> /4%		GOOGI	MFTA	-0.35%	COST -1.01%	LLY -0.41%	ABBV -0.48%	MRK -0.42%	UI -0.1	<b>NH</b> 95%
MSFT -1.70%	ADBE -1.04%		AMD -1.46%	INTC -1.34%	QCOM -0.64%	-1.94%	-0.62%	<b>PG</b>	<b>KO</b> -0.82%	INI	PFE -1.76%	BMY -1.05%	ELV -1.29%	CI -1.99% HUM
	PANW SNPS +0.12% -0.78%	AVGO -2.12%	TXN	ADI -1.76%	NXPI -1.58%			KVUE EL		-0.77% MEDICAL DEVICES	AMGN +0.82%	GILD +0.94%	CVS -1.71%	-1.61% CNC -0.82% IOTECHN
	FTNT -1.17 FLT GEN	SOFTWARE - A	-183%	-0.93%	ON 456 TION TECH	CMCSA VZ -1.40% -0.51% NF			-1.04% KDP -0.98	ABT SYK -1.27% -4.81%	<b>TMO</b> -0.81%	ISRG -0.86%	BDX -2.19% ¥0	RTX REGN 1.29% -0.27
		CRM -1.87%	INTU -1.01%	ACN -2.01%	IBM -2.21%	TMUS +0.32% -0.27% -0.27% -0.27%	** -1.19% EA TTWO		ADM BEVERAG G STZ	MDT EW -1.27% +1.31	DHR -1.90%	RMD BAX MEDICA	M DRUG	IRNA MEDICAL
		NOW -2.65%	S ADSK -1.95% FICO TYL	L FI -1.32	H IT CDW 5 -2.21 -2.32	CONSUMER CYCLICAL	UTO MANUFACTUR	MO -0.45% +1.23%	HSY SYY KR	BSX -1.80%	-1.10 -1.08	-0.76 CAH	ZTS -0.48%	HCA GEHC
-2.34%		-1.859 COMMUNICA -1.44% MSI -1.14%	* PTC SEMICON AMA -1.84% LRCX KI -2.35%	ANET	J ELECTR T APH TEL -2.04 "IFIC FSLR TDY	<b>AMZN</b> -4.03%	<b>TSLA</b> -1.16%	SPECIALTY INDUSTRI GE ETN -1.60% -2.24% ITW -1.95% -2.12 or -1.87% -1.95% -2.12 or -1.87% -1.06 IR 1.	AEROSPACE & I BA L -1.58% -1 IIS RTX GD +0.29% -0.52	DEFE REIT-SPECIAL MT AMT CI -1.70 CCI TDG REIT-INDU I -1.508 PLD PSA	DLR WY SBAC REIT- ESS	OM .15%	D CVX -0.54%	OIL & GA MPC PSX -2.06 -0.70 VLO OIL & GA SLB
FINANCIAL CREDIT SERVICES BANKS - DIVERSIFIED		FINANCIAL	DATA & ST	CAPITAL	MARKETS			-255% ROK -1.89 IEX	NOC LHX +0.24%	EXR	CPT OIL & G	AS E&P		HAL BKR -0.46 -0.53
<b>V</b> JPM	BAC -1.56%	SPGI -1.99%	ICE MCC -1.79 -2.29	<sup>9</sup> MS -1.04% 9	GS 5CHW -2.08%	HOME IMPROVEMENT SPECIALTY I AZO -1.04	F GM	CAT -1.73% DCAD	HON -2.07% -1.10	P SPG VTR	CBRE -0.77%	EOG PX +0.98% +10 OXY HE -0.02% +0.0	S EQT	OIL&GA WMB KMI
AXP	C -1.95%	CB PGR	E - PROP	INSURANC MMC	BANKS - R USB PNC -1.70 -0.84	-1.29% LOW TSCO -1.86% FOOTWE NKE	-1.21% RESIDEN LODGI DHI MAR	RAILROADS	UILDING P SPEC	GPN NEE	GULATED ELE	SPECIAL SPECIAL	MATERIA TY CHE	ALS AGRICUL CTVA
-1.77% PYPL +0.12% COS DFS COS SYF	(- <b>B</b>	IG ASSET MAN BLK -1.65% BH -1.4	L NAGEME IP I I I I I I I I I I I I I I I I I I	AON -1.25% AJG INSURANCE MET AFL	TFC RF -2.05 MTB LIFE EG -1.98	SBUX -0.75%         -0.47%           MCD -0.75%         -0.80%           CMG -0.65%         YUM -0.45%           DRI         -0.44%	LEN 4001 PACKAGI AUTO P IP RESORTS LVS	INTEGRATED FREI	WM RSG LUV 0.94 RSG LUV INDUST J TRUCK	VRSK 0.53% EFX URI -3.15% SNA DUK -3.15% DUK -3.06 	XEL         ED           -2.98         -2.66           ES         DTE           FE         PPL           ETR         LNT	EG 331 SHW -0.91% WK APD -1.00%	ECL DD 1.86 DD LYB ALB PPG IFF	-2.26 FCX -1.69 STLD BUILDI CHEMI
https://finviz.com	n/map	.ashx	ad via			aru aranh				-3% -2%	-1% 0%	5 +1%	+2%	+3%

# ENCODINGS

market cap (Q)  $\rightarrow$  rectangle size mkt sector (N), mkt cap (Q)  $\rightarrow$  rect. pos. loss vs. gain (N, O)  $\rightarrow$  color hue magnitude of loss or gain (Q)  $\rightarrow$  color value

















# **COMBINATORICS OF ENCODINGS**

### **Challenge:**

Assume **k** visual attributes/channels and **n** data fields Pick the best encoding from the exponential number of possibilities  $(n+1)^k$ 

### PRINCIPLES

### Challenge

Assume **k** visual attributes/channels and **n** data fields Pick the best encoding from the exponential number of possibilities  $(n+1)^k$ 

### **Principle of Consistency**

Properties of image (visual variables) should match properties of data

### **Principle of Importance Ordering**

Encode most important information in the most effective way

# **EXPRESSIVENESS CRITERIA** [Mackinlay 1986]

### **Expressiveness**

A set of facts is expressible in a visual language if the sentences (i.e., the visualizations) in the language express *all* the facts in the set of data, and *only* the facts in the data.

# **CANNOT EXPRESS ALL THE FACTS**

### Horizontal dot plot

A one-to-many (1  $\rightarrow$  N) relation cannot be expressed in a single horizontal dot plot because multiple tuples are mapped to the same position

		000																
		0	5	10	15	20	25	30	35	40	45	50	55	60	65	70	75	80
										Value	•							
	petal	_		0000	0000	•												
I. Setosa	sepal					-					00000			×				
L. Vorginico	petal										•				•••	• •		
1. verginica	sepal											٠	~			•••• ••	••••	• •
I. Versicolor	petal							٠	0 000	00000	00000		• •					
1. Versicolor	sepal											0000		00000				
		0		10		20		 30		40		50		60		 70		 80
										Valu	е							

![](_page_29_Figure_1.jpeg)

# **EFFECTIVENESS CRITERIA** [Mackinlay 1986]

### Effectiveness

A visualization is more effective than another visualization if the information conveyed by one visualization is more readily *perceived* than the information in the other visualization.

Subject of the Perception Lecture

![](_page_30_Figure_1.jpeg)

# AUTOMATIC CHART DESIGN [Mackinlay 1986]

APT - "A Presentation Tool"

### User formally specifies data model and type

Input: list of data variables ordered by importance

### APT searches over the design space

Tests expressiveness of each visual encoding (rule-based) Generates encodings that pass test Rank by perceptual effectiveness criteria

### Outputs most effective visualization

![](_page_31_Picture_1.jpeg)

![](_page_31_Figure_3.jpeg)

### LIMITATIONS

Does not cover many visualization techniques

Networks, maps, diagrams Also, 3D, animation, illustration, ...

Does not consider interaction Does not consider semantics or conventions Assumes single visualization as output

#### 82

### **SUMMARY**

#### **Formal specification**

Data model: tidy data, N,O,Q types Image model: marks, visual attributes/channels Encodings map data to mark attributes/channels

### Choose expressive and effective encodings

Rule-based test of expressiveness Perceptual effectiveness rankings