

G53FIV: Fundamentals of Information Visualization

Lecture 7: Visualization with R – Advanced and Visualization Tools

Ke Zhou
School of Computer Science
Ke.Zhou@nottingham.ac.uk

<https://moodle.nottingham.ac.uk/course/view.php?id=96914>

Last Lecture

Visualization with R

R is a tool for...

Data Manipulation

- connecting to data sources
- slicing & dicing data

Modeling & Computation

- statistical modeling
- numerical simulation

Data Visualization

- visualizing fit of models
- composing statistical graphics

munge



model



visualize

Building a Plot in **ggplot2**

data to visualize (a data frame)

map variables to **aes**thetic attributes

geometric objects – what you see (points, bars, etc)

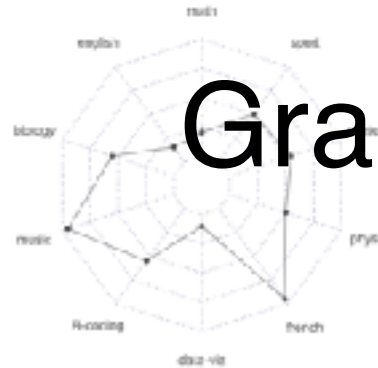
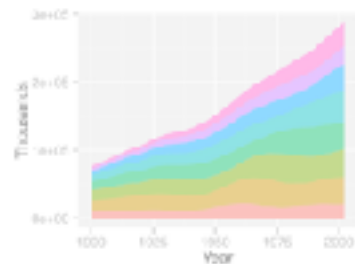
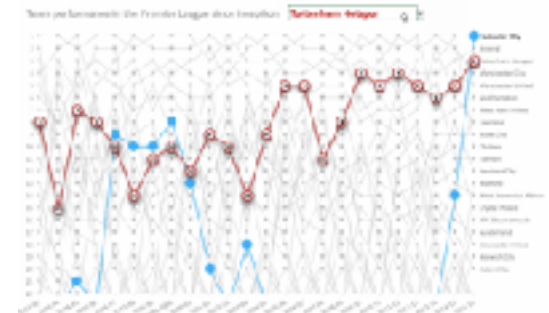
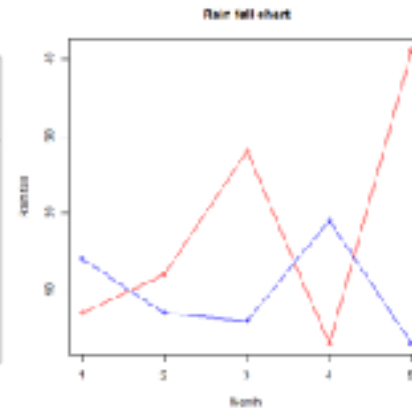
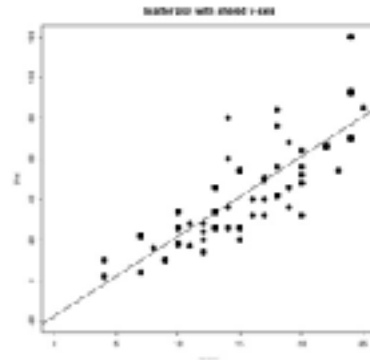
scales map values from data to aesthetic space

```
ggplot(iris) + geom_point(aes(x = Sepal.Length, y = Sepal.Width))
```

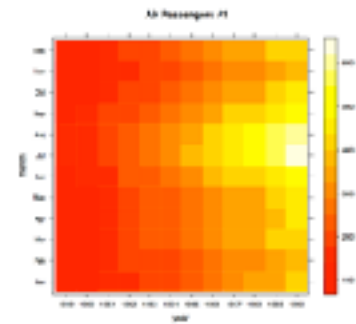
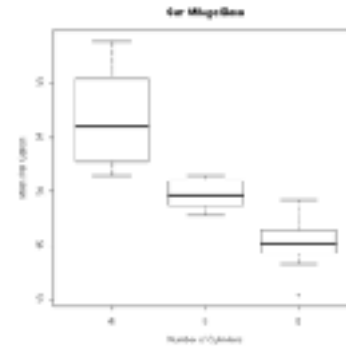
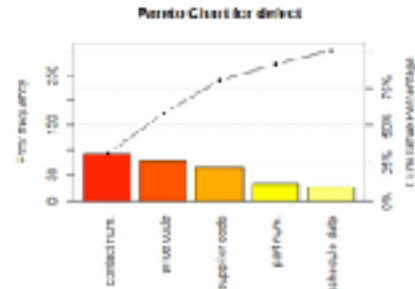
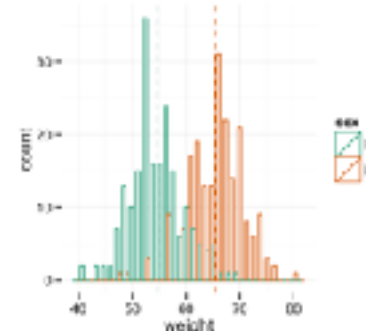
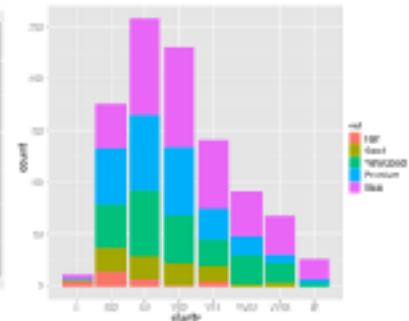
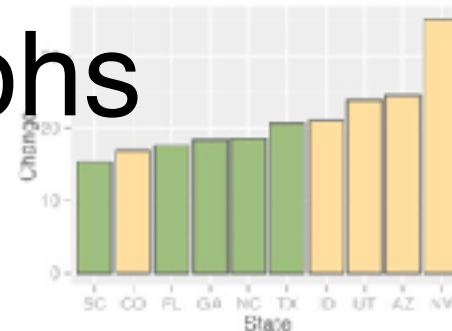
Data

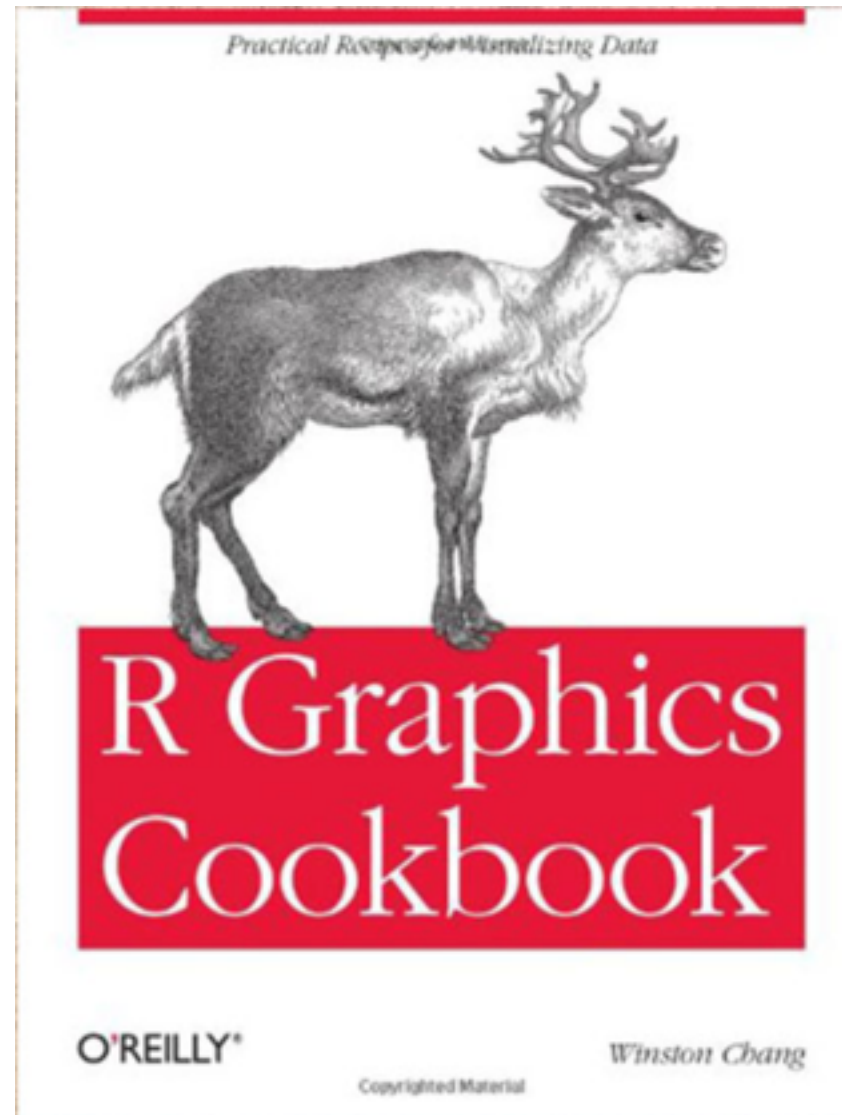
Geometric objects to display

Aesthetics map variables to scales

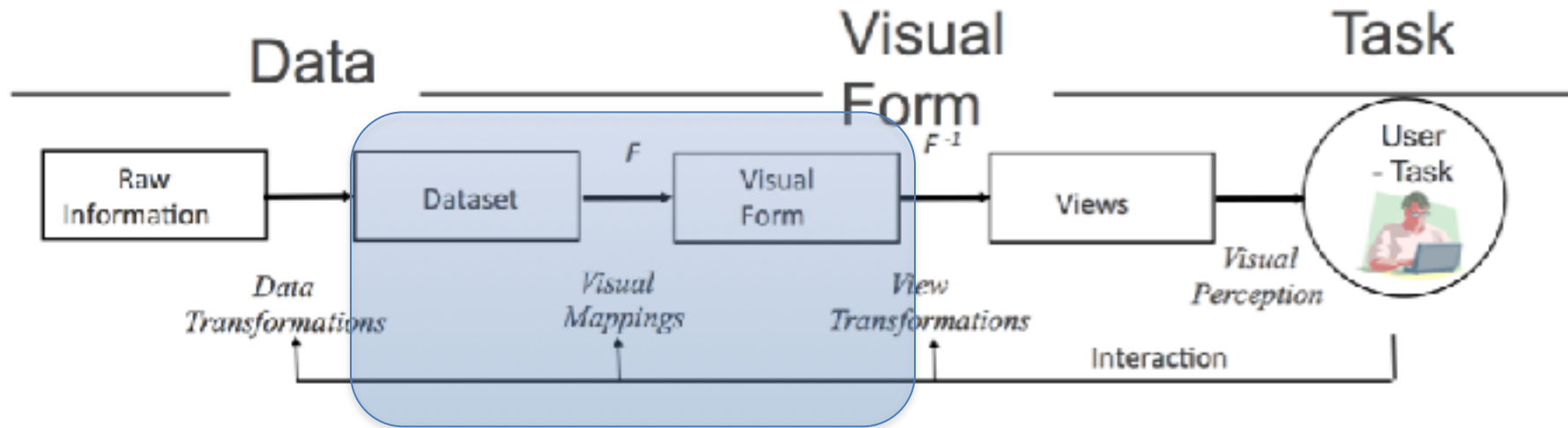


Graphs

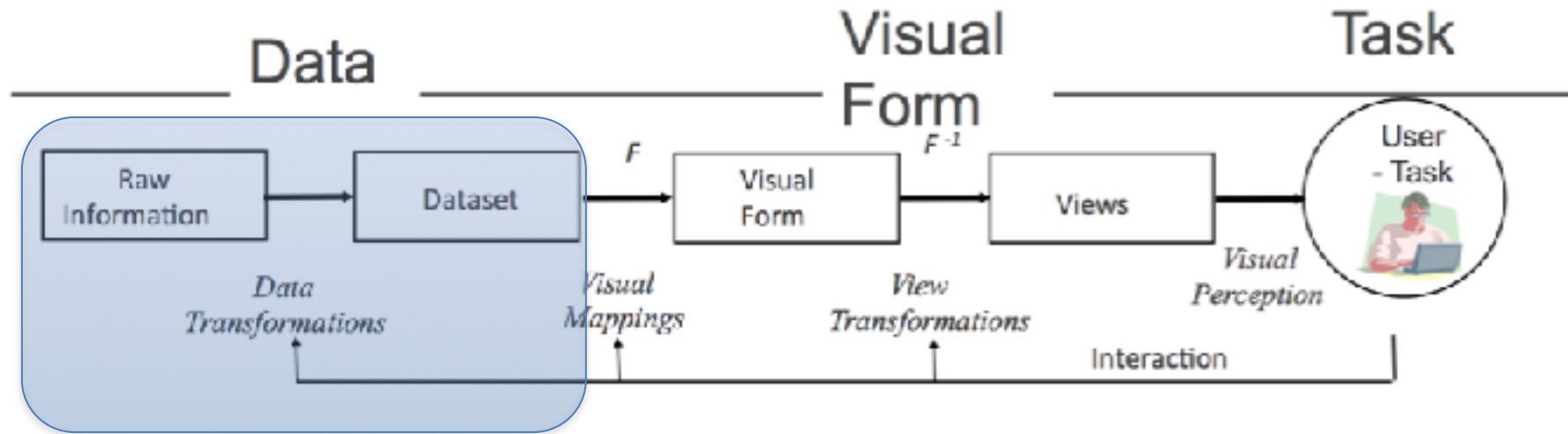




Seven Stages of Visualization



Seven Stages of Visualization



G53FIV Early Module Feedback

- Survey (on Moodle)
 - Anonymous
 - Seeking for constructive feedback
 - Two compulsory multiple questions
 - A few open-ended questions
 - If you have specific feedback or suggestions for improvement.
- A summary of the survey (and action points for improvements) will be presented.

Overview

- Data Manipulations with R
- Visualization Tools

Data Manipulations with R

Transform Data: A Swiss-Army Knife

- Indexing
- Three ways to index into a data frame
 - Array of integer indices
 - Array of character names
 - Array of logical Booleans
- Examples:
 - `df[1:3,]`
 - `df[c("New York", "Chicago"),]`
 - `df[c(TRUE, FALSE, TRUE, TRUE),]`

	A	B	C	D
1	year	sex	smoke	city
2	1992	M	Y	1
3	1992	M	Y	2
4	1992	M	Y	2
5	1992	M	Y	2
6	1992	M	Y	2
7	1992	M	Y	2
8	1992	M	Y	2
9	1992	M	Y	2
10	1992	M	Y	2
11	1992	M	Y	2
12	1992	M	Y	2
13	1992	M	Y	2
14	1992	M	Y	2
15	1992	M	Y	2
16	1992	M	Y	2
17	1992	M	Y	2
18	1992	M	Y	2
19	1992	M	Y	2
20	1992	M	Y	2
21	1992	M	Y	2
22	1992	M	Y	2
23	1992	M	Y	2



Transform Data: A Swiss-Army Knife

- **subset** – extract subsets meeting some criteria

```
subset(Insurance, District==1)
```

```
subset(Insurance, Claims < 20)
```

	A	B	C	D
1	year	age	sex	inc
2	1990	5	M	1
3	1990	5	M	2
4	1990	9	M	1
5	1990	1	M	2
6	1990	10	M	1
7	1990	10	M	2
8	1990	10	M	1
9	1990	10	M	2
10	1990	15	M	1
11	1990	20	M	2
12	1990	11	M	1
13	1990	15	M	2
14	1990	20	M	1
15	1990	20	M	2
16	1990	20	M	1
17	1990	25	M	2
18	1990	40	M	1
19	1990	40	M	2
20	1990	45	M	1
21	1990	48	M	2
22	1990	80	M	1
23	1990	50	M	2
24	1990	65	M	1



Transform Data: A Swiss-Army Knife

- **subset** – extract subsets meeting some criteria
`subset(Insurance, District==1)`
`subset(Insurance, Claims < 20)`
- **transform** – add or alter a column of a data frame
`transform(Insurance, Propensity=Claims/Holders)`

	A	B	C	D
1	1990	100	1000	100
2	1990	0	0	0
3	1990	0	0	0
4	1990	1	0	0
5	1990	5	0	0
6	1990	10	0	0
7	1990	10	0	0
8	1990	10	0	0
9	1990	10	0	0
10	1990	10	0	0
11	1990	10	0	0
12	1990	10	0	0
13	1990	10	0	0
14	1990	10	0	0
15	1990	10	0	0
16	1990	10	0	0
17	1990	10	0	0
18	1990	10	0	0
19	1990	10	0	0
20	1990	10	0	0
21	1990	10	0	0
22	1990	10	0	0
23	1990	10	0	0
24	1990	10	0	0



Transform Data: A Swiss-Army Knife

- **subset** – extract subsets meeting some criteria
`subset(Insurance, District==1)`
`subset(Insurance, Claims < 20)`
- **transform** – add or alter a column of a data frame
`transform(Insurance, Propensity=Claims/Holders)`
- **cut** – cut a continuous value into groups
`cut(Insurance$Claims, breaks=c(-1,100,Inf), labels=c('lo','hi'))`

	A	B	C	D
1	year	age	sex	region
2	1992	5	D	1
3	1992	5	D	2
4	1992	8	D	3
5	1992	5	D	4
6	1992	32	D	5
7	1992	32	D	6
8	1992	15	D	7
9	1992	12	D	8
10	1992	35	D	9
11	1992	32	D	10
12	1992	31	D	11
13	1992	35	D	12
14	1992	35	D	13
15	1992	35	D	14
16	1992	32	D	15
17	1992	32	D	16
18	1992	32	D	17
19	1992	32	D	18
20	1992	32	D	19
21	1992	32	D	20
22	1992	32	D	21
23	1992	32	D	22
24	1992	32	D	23
25	1992	32	D	24
26	1992	32	D	25
27	1992	32	D	26
28	1992	32	D	27
29	1992	32	D	28
30	1992	32	D	29
31	1992	32	D	30
32	1992	32	D	31
33	1992	32	D	32
34	1992	32	D	33
35	1992	32	D	34
36	1992	32	D	35
37	1992	32	D	36
38	1992	32	D	37
39	1992	32	D	38
40	1992	32	D	39
41	1992	32	D	40
42	1992	32	D	41
43	1992	32	D	42
44	1992	32	D	43
45	1992	32	D	44
46	1992	32	D	45
47	1992	32	D	46
48	1992	32	D	47
49	1992	32	D	48
50	1992	32	D	49
51	1992	32	D	50
52	1992	32	D	51
53	1992	32	D	52
54	1992	32	D	53
55	1992	32	D	54
56	1992	32	D	55
57	1992	32	D	56
58	1992	32	D	57
59	1992	32	D	58
60	1992	32	D	59
61	1992	32	D	60
62	1992	32	D	61
63	1992	32	D	62
64	1992	32	D	63
65	1992	32	D	64
66	1992	32	D	65
67	1992	32	D	66
68	1992	32	D	67
69	1992	32	D	68
70	1992	32	D	69
71	1992	32	D	70
72	1992	32	D	71
73	1992	32	D	72
74	1992	32	D	73
75	1992	32	D	74
76	1992	32	D	75
77	1992	32	D	76
78	1992	32	D	77
79	1992	32	D	78
80	1992	32	D	79
81	1992	32	D	80
82	1992	32	D	81
83	1992	32	D	82
84	1992	32	D	83
85	1992	32	D	84
86	1992	32	D	85
87	1992	32	D	86
88	1992	32	D	87
89	1992	32	D	88
90	1992	32	D	89
91	1992	32	D	90
92	1992	32	D	91
93	1992	32	D	92
94	1992	32	D	93
95	1992	32	D	94
96	1992	32	D	95
97	1992	32	D	96
98	1992	32	D	97
99	1992	32	D	98
100	1992	32	D	99



Transform Data: A Swiss-Army Knife

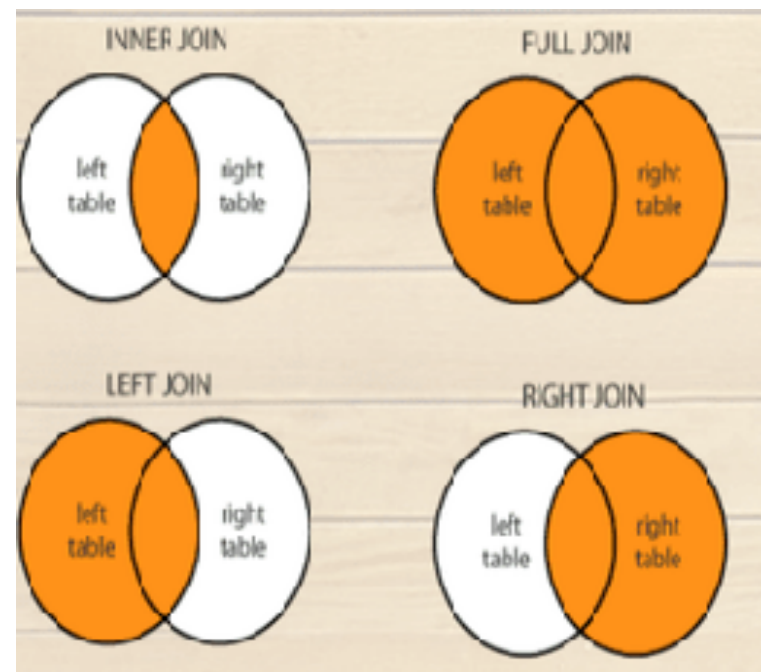
- **subset** – extract subsets meeting some criteria
`subset(Insurance, District==1)`
`subset(Insurance, Claims < 20)`
- **transform** – add or alter a column of a data frame
`transform(Insurance, Propensity=Claims/Holders)`
- **cut** – cut a continuous value into groups
`cut(Insurance$Claims, breaks=c(-1,100,Inf), labels=c('lo','hi'))`
- Put it all together: create a new, transformed data frame

```
transform(subset(Insurance, District==1),  
  ClaimLevel=cut(Claims, breaks=c(-1,100,Inf),  
    labels=c('lo','hi')))
```



Joining Two Data Frames

- `inner_join(df1, df2, by = "common_column")`
- `?join`
 - `Left_join`, `right_join`
 - `Inner_join`, `outer_join`
- `merge(x=df1, y=df2, by.x="id", by.y="bid")`



Libraries for Data Manipulations

- Packages
 - plyr
 - data.table
 - reshape2
 - doBY
 - sqldf
 - and many more

dplyr: A Grammar of Data Manipulation

- Very intuitive, once you understand the basics
- Very fast
 - Created with execution times in mind
- Easy for those migrating from the SQL world
- When written well, your code reads like a “recipe”
- “Code the way you think”

<https://cran.rstudio.com/web/packages/dplyr/vignettes/introduction.html>

Pipe Operator

- Library(maggritr)
 - A R package launched on Jan 2014
 - A “magic” operator called the PIPE was introduced
 - %>%
 - i.e. “AND THEN”, “PIPE TO”

```
round(sqrt(1000), 3)
```

```
library(maggritr)
```

```
1000 %>% sqrt %>% round()
```

```
1000 %>% sqrt %>% round(., 3)
```

Take 1000, and then its sqrt
And then round it



dplyr

- dplyr takes the `%>%` operator and uses it to great effect for manipulating data frames
 - Works only with data frames
 - 5 basic “verbs” work for 90% of data

Verbs	What does it do?
<code>filter()</code>	Select a subset of ROWS by conditions
<code>arrange()</code>	Reorders ROWS in a data frame
<code>select()</code>	Select the COLUMNS of interest
<code>mutate()</code>	Create new columns based on existing columns (mutations!)
<code>summarise()</code>	Aggregate values for each group, reduces to single value

5 Basic Verbs

- **FILTE_RRows**



- **SELE_{CT} Column Types**



- **Ar_Range Rows (SORT)**



- **Mutate (into something new)**



- **Summarize by Groups**



Movies dataset

title	year	budget	votes	length	Docume ntary	rating	...
Titanic	1997	200,000,000	1000	195	0	7.8	...
Leon	1994	16,000,000	500	90	0	8.6	...
McQueen	2018	52,000,000	200	91	1	7.9	...

Filter()



- Usage: **`filter(data, condition)`**
 - Returns a subset of rows
 - Multiple conditions can be supplied.
 - They are combined with an AND

```
movies_with_budgets <- filter(movies_df, !is.na(budget))  
filter(movies, Documentary==1)  
filter(movies, Documentary==1) %>% nrow()
```


Filter()



- Usage: **`filter(data, condition)`**
 - Returns a subset of rows
 - Multiple conditions can be supplied.
 - They are combined with an AND

```
movies_with_budgets <- filter(movies_df, !is.na(budget))  
filter(movies, Documentary==1)  
filter(movies, Documentary==1) %>% nrow()  
good_comedies <- filter(movies, rating > 9, Comedy==1)  
dim(good_comedies) #171 movies
```

Filter()



- Usage: **filter(data, condition)**
 - Returns a subset of rows
 - Multiple conditions can be supplied.
 - They are combined with an AND

```
movies_with_budgets <- filter(movies_df, !is.na(budget))
filter(movies, Documentary==1)
filter(movies, Documentary==1) %>% nrow()
good_comedies <- filter(movies, rating > 9, Comedy==1)
dim(good_comedies) #171 movies
```

```
#' Let us say we only want highly rated comedies, which a lot
of people have watched, made after year 2000.
```

```
movies %>%
  filter(rating > 8, Comedy==1, votes > 100, year > 2000)
```

Select()

- Usage:

`select(data, columns)`

```
movies_df <- tbl_df(movies)
select(movies_df, title, year, rating) #Just the columns we want to see
select(movies_df, -c(r1:r10)) #we don't want certain columns
```



Select()

- Usage:

`select(data, columns)`

```
movies_df <- tbl_df(movies)
select(movies_df, title, year, rating) #Just the columns we want to see
select(movies_df, -c(r1:r10)) #we don't want certain columns

#You can also select a range of columns from start:end
select(movies_df, title:votes) # All the columns from title to votes
select(movies_df, -c(budget, r1:r10, Animation, Documentary, Short, Romance))
```



Select()

- Usage:

`select(data, columns)`

```
movies_df <- tbl_df(movies)
select(movies_df, title, year, rating) #Just the columns we want to see
select(movies_df, -c(r1:r10)) #we don't want certain columns

#You can also select a range of columns from start:end
select(movies_df, title:votes) # All the columns from title to votes
select(movies_df, -c(budget, r1:r10, Animation, Documentary, Short, Romance))

select(movies_df, contains("r")) # Any column that contains 'r' in its name
select(movies_df, ends_with("t")) # All vars ending with "t"
```



Select()

- Usage:

`select(data, columns)`



```
movies_df <- tbl_df(movies)
select(movies_df, title, year, rating) #Just the columns we want to see
select(movies_df, -c(r1:r10)) #we don't want certain columns

#You can also select a range of columns from start:end
select(movies_df, title:votes) # All the columns from title to votes
select(movies_df, -c(budget, r1:r10, Animation, Documentary, Short, Romance))

select(movies_df, contains("r")) # Any column that contains 'r' in its name
select(movies_df, ends_with("t")) # All vars ending with "t"

select(movies_df, starts_with("r")) # Gets all vars staring with "r"
#The above is not quite what we want. We don't want the Romance column
select(movies_df, matches("r[0-9]")) # Columns that match a regex.
```

Arrange()



Usage:

`arrange(data, column_to_sort_by)`

- Returns a reordered set of rows
- Multiple inputs are arranged from left-to-right

```
movies_df <- tbl_df(movies)
arrange(movies_df, rating) #but this is not what we want
arrange(movies_df, desc(rating))
```


Arrange()



Usage:

`arrange(data, column_to_sort_by)`

- Returns a reordered set of rows
- Multiple inputs are arranged from left-to-right

```
movies_df <- tbl_df(movies)
arrange(movies_df, rating) #but this is not what we want
arrange(movies_df, desc(rating))
#Show the highest ratings first and the latest year...
#Sort by Decreasing Rating and Year
arrange(movies_df, desc(rating), desc(year))
```


Arrange()



Usage: `arrange(data, column_to_sort_by)`

- Returns a reordered set of rows
- Multiple inputs are arranged from left-to-right

```
movies_df <- tbl_df(movies)
arrange(movies_df, rating) #but this is not what we want
arrange(movies_df, desc(rating))
#Show the highest ratings first and the latest year...
#Sort by Decreasing Rating and Year
arrange(movies_df, desc(rating), desc(year))
```

What's the difference between these two?

```
arrange(movies_df, desc(rating), desc(year))
arrange(movies_df, desc(year), desc(rating))
```

Mutate()



- Usage:

```
mutate(data, new_col = func(oldcolumns))
```

- Creates new columns, that are functions of existing variables

```
movies_with_budgets <- filter(movies_df, !is.na(budget))  
mutate(movies_with_budgets, costPerMinute = budget/length) %>%  
  select(title, costPerMinute)
```

Group_by() and Summarize()

```
group_by(data, column_to_group) %>%  
  summarize(function_of_variable)
```

- Group_by creates groups of data
- Summarize aggregates the data for each group

```
by_rating <- group_by(movies_df, rating)
```

```
by_rating %>% summarize(n())
```

Group_by() and Summarize()

```
group_by(data, column_to_group) %>%  
  summarize(function_of_variable)
```

- Group_by creates groups of data
- Summarize aggregates the data for each group

```
by_rating <- group_by(movies_df, rating)
```

```
by_rating %>% summarize(n())
```

```
avg_rating_by_year <-  
  group_by(movies_df, year) %>%  
  summarize(avg_rating = mean(rating))
```

Chain the “Verbs” Together

- Chain them together

```
producers_nightmare <-  
  filter(movies_df, !is.na(budget)) %>%  
  mutate(costPerMinute = budget/length) %>%  
  arrange(desc(costPerMinute)) %>%  
  select(title, costPerMinute)
```

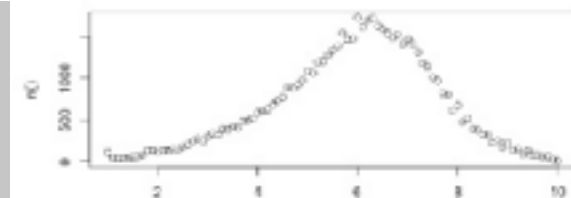
Chain the “Verbs” Together

- Chain them together

```
producers_nightmare <-  
  filter(movies_df, !is.na(budget)) %>%  
  mutate(costPerMinute = budget/length) %>%  
  arrange(desc(costPerMinute)) %>%  
  select(title, costPerMinute)
```

- Can also be fed to a “plot” command

```
movies %>%  
  group_by(rating) %>%  
  summarize(n()) %>%  
  plot() # plots the histogram of movies by Each value of rating
```



Practice

- Find all the post-2000 comedy movies with budget of over \$1,000,000, rank them by rating in the decreasing order, and output their title and rating

Practice

- Find all the post-2000 comedy movies with budget of over \$1,000,000, rank them by rating in the decreasing order, and output their title and rating
- **`comedies <- filter(movies_df, year > 2000, Comedy = 1, budget > 1000000) %>%`**

Practice

- Find all the post-2000 comedy movies with budget of over \$1,000,000, rank them by rating in the decreasing order, and output their title and rating
- **`comedies <- filter(movies_df, year > 2000, Comedy = 1, budget > 1000000) %>%`**
- **`arrange(desc(rating)) %>%`**
- **`select(title, rating)`**

Visualization Tools

Visualization Tools

Chart Typologies

Excel, Many Eyes, Google Charts

Visual Analysis Grammars

VizQL, ggplot2

Visualization Grammars

Protovis, D3.js

Component Architectures

Prefuse, Flare, Improvise, VTK

Graphics APIs

Processing, OpenGL, Java2D

Visualization Tools

Chart Typologies

Excel, Many Eyes, Google Charts

Charting
Tools

Visual Analysis Grammars

VizQL, ggplot2

Declarative
Languages

Visualization Grammars

Protovis, D3.js

Component Architectures

Prefuse, Flare, Improvise, VTK

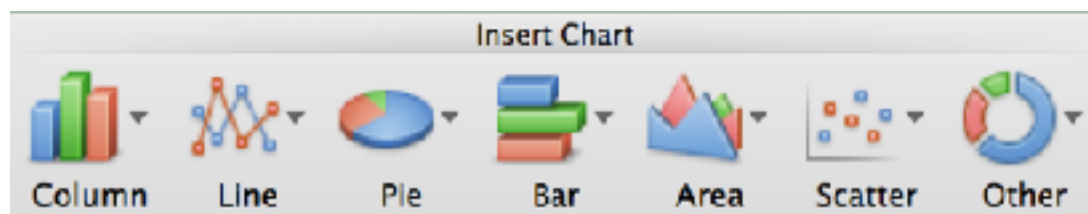
Programming
Toolkits

Graphics APIs

Processing, OpenGL, Java2D

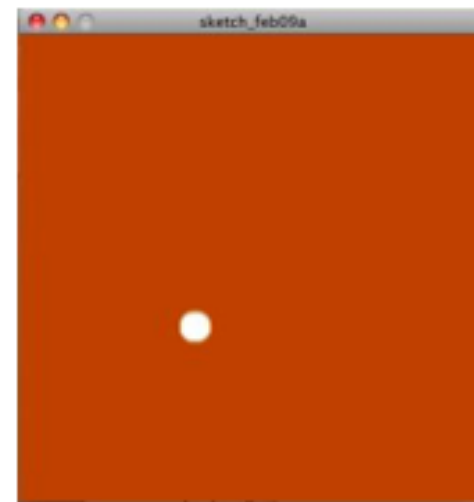
Chart Typology (Charting Tools)

- Pick from a stock of templates
- Easy-to-use but limited expressiveness
- Prohibits novel designs, new data types



Graphics APIs (Programming Toolkits)

- Processing.org
 - Java based
 - not specifically designed for InfoVis
 - Well documented, lots of tutorials (even books)



Graphics APIs can be very powerful

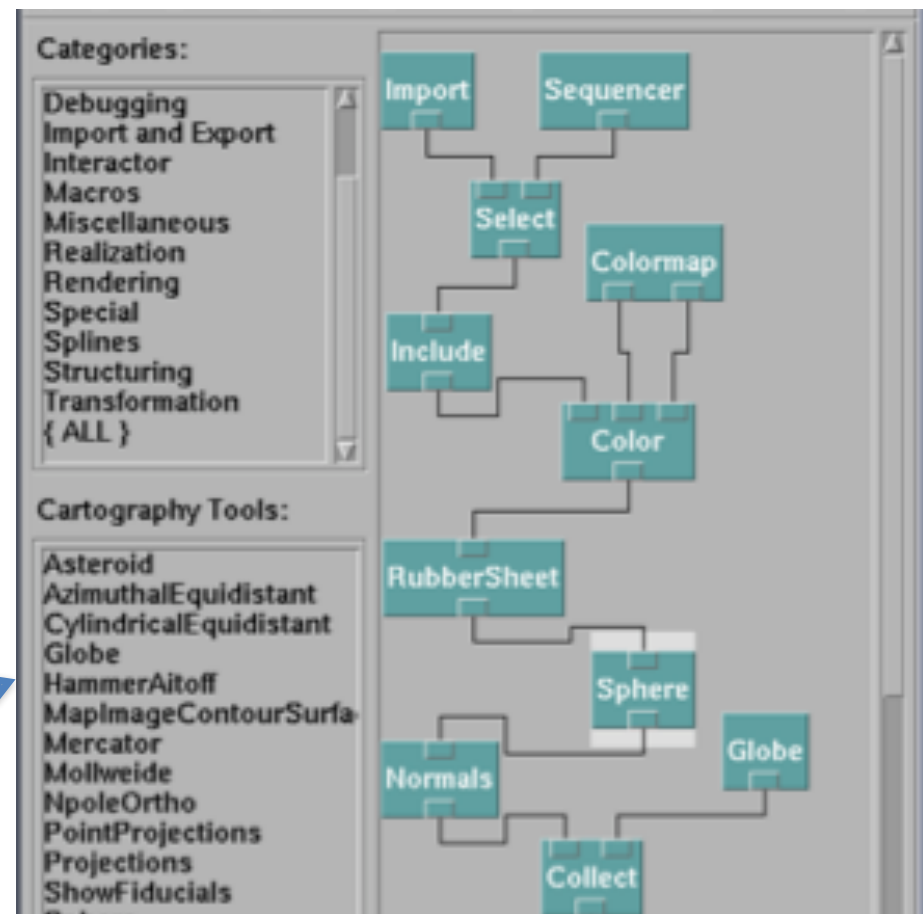
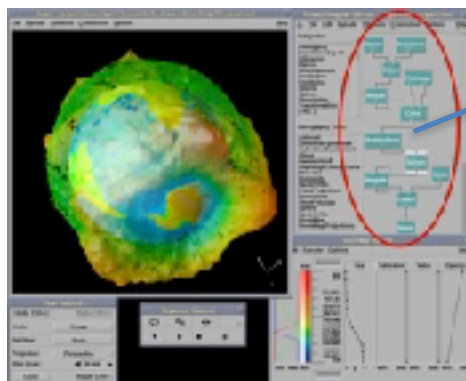


US Air Traffic Visualization

Dr. Ke Zhou (<http://www.cs.nott.ac.uk/~pszkhz/>)

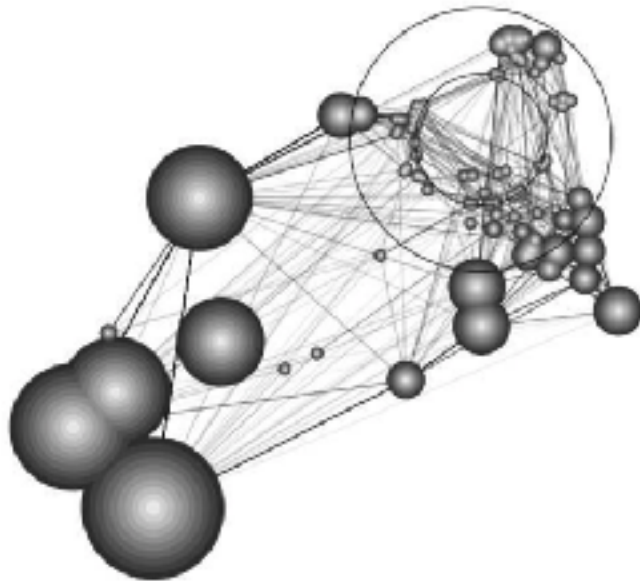
Component Architectures (Programming Toolkits)

- Permits more combinatorial possibilities
- Novel views require new operators, which requires software engineering.

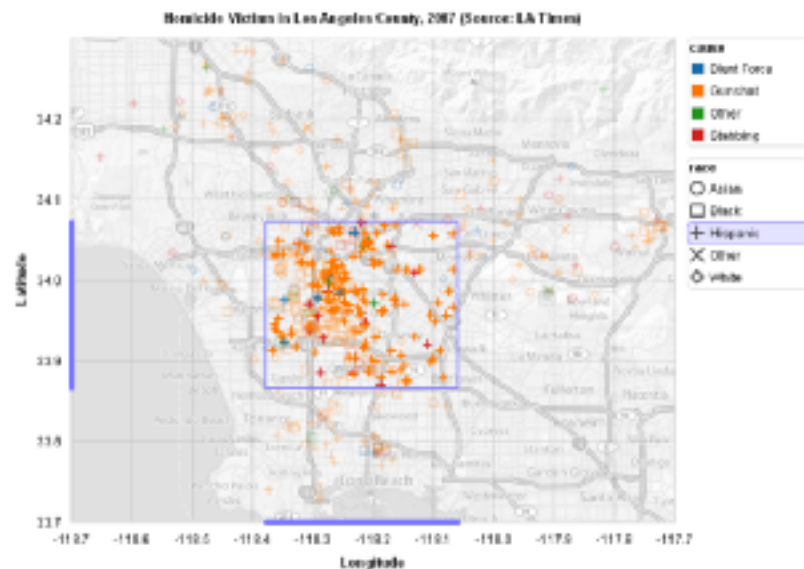


Prefuse & Flare

- Operator-based toolkits for visualization design
 $\text{Vis} = (\text{Input Data} \rightarrow \text{Visual Objects}) + \text{Operators}$



Prefuse (<http://prefuse.org>)



Flare (<http://flare.prefuse.org>)

Comparison



- **Chart Typology**
 - Pick from a stock of templates
 - Easy-to-use but limited expressiveness
 - Prohibits novel designs, new data types

- **Component Architecture**
 - Permits more combinatorial possibilities
 - Novel views require new operators, which requires software engineering.

The Grammar of Graphics (Declarative Languages)

- Programming by describing what, not how
- Separate specification (what you want) from execution (how it should be computed)
- In contrast to imperative programming, where you must give explicit steps.

```
d3.selectAll("rect")  
  .data(my_data)  
  .enter().append("rect")  
  .attr("x", function(d) { return xscale(d.foo); })  
  .attr("y", function(d) { return yscale(d.bar); })
```

Building a Plot in ggplot2

data to visualize (a data frame)

map variables to **aes**thetic attributes

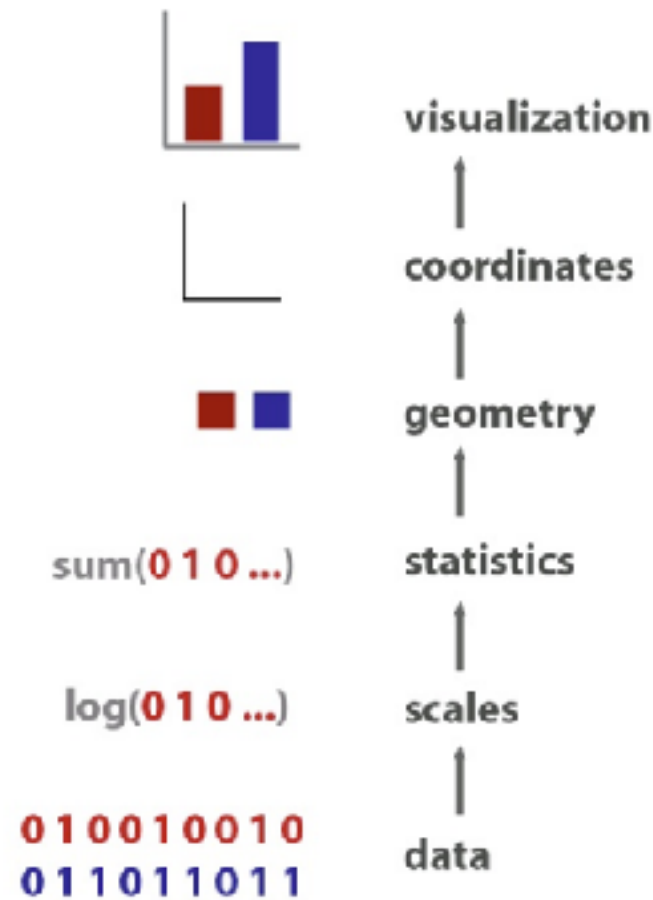
geometric objects – what you see
(points, bars, etc)

scales map values from data to
aesthetic space

faceting subsets the data to show
multiple plots

statistical transformations – summarize
data

coordinate systems put data on plane of
graphic



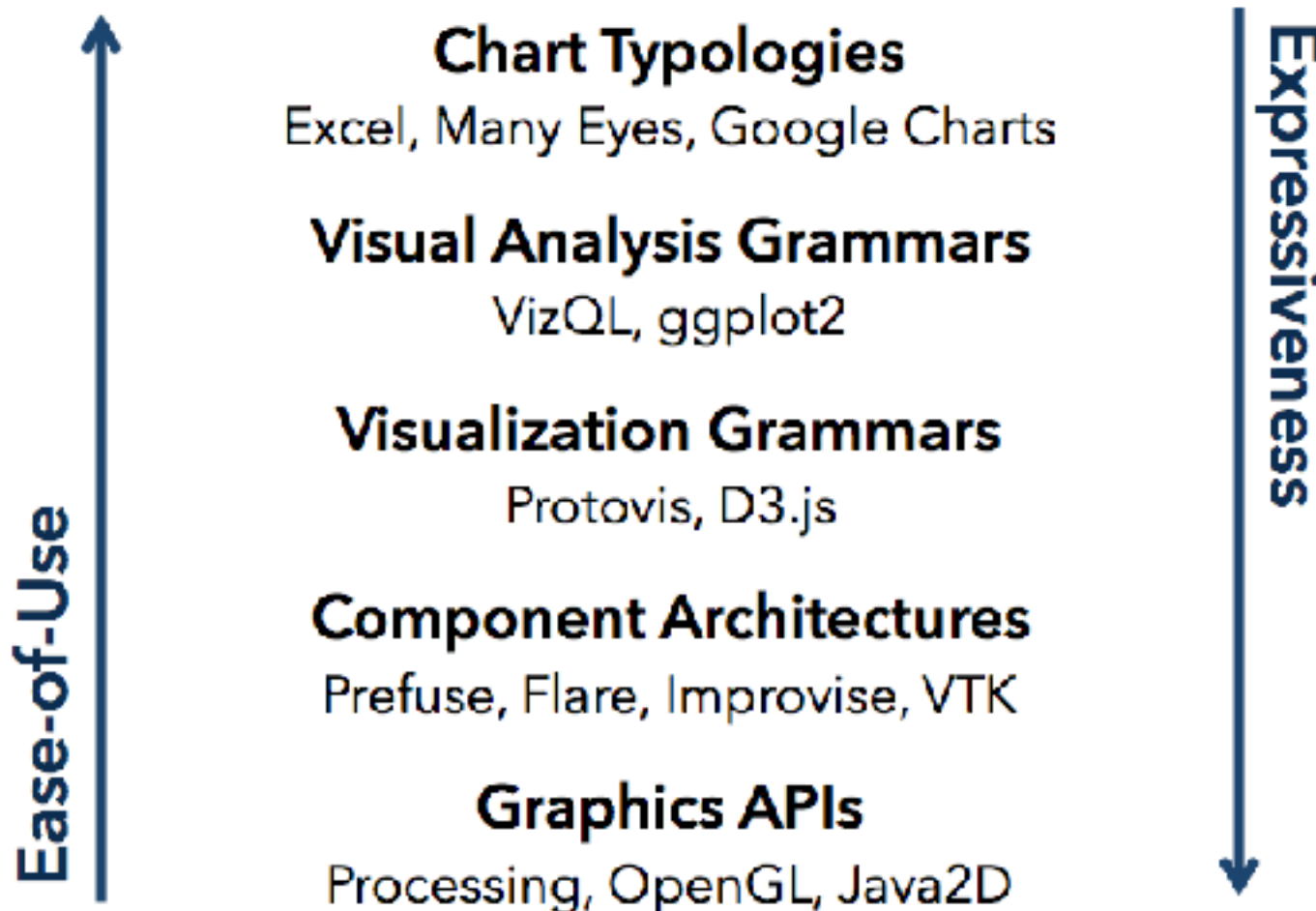
The Advantages of Declarative Languages

- **Faster iteration.** Less code. Larger user base.
- **Better visualization.** Smart defaults.
- **Reuse.** Write-once, then re-apply.
- **Performance.** Optimization, scalability.
- **Portability.** Multiple devices, renderers, inputs.
- **Programmatic generation.** Write programs which output visualizations. Automated search & recommendation.

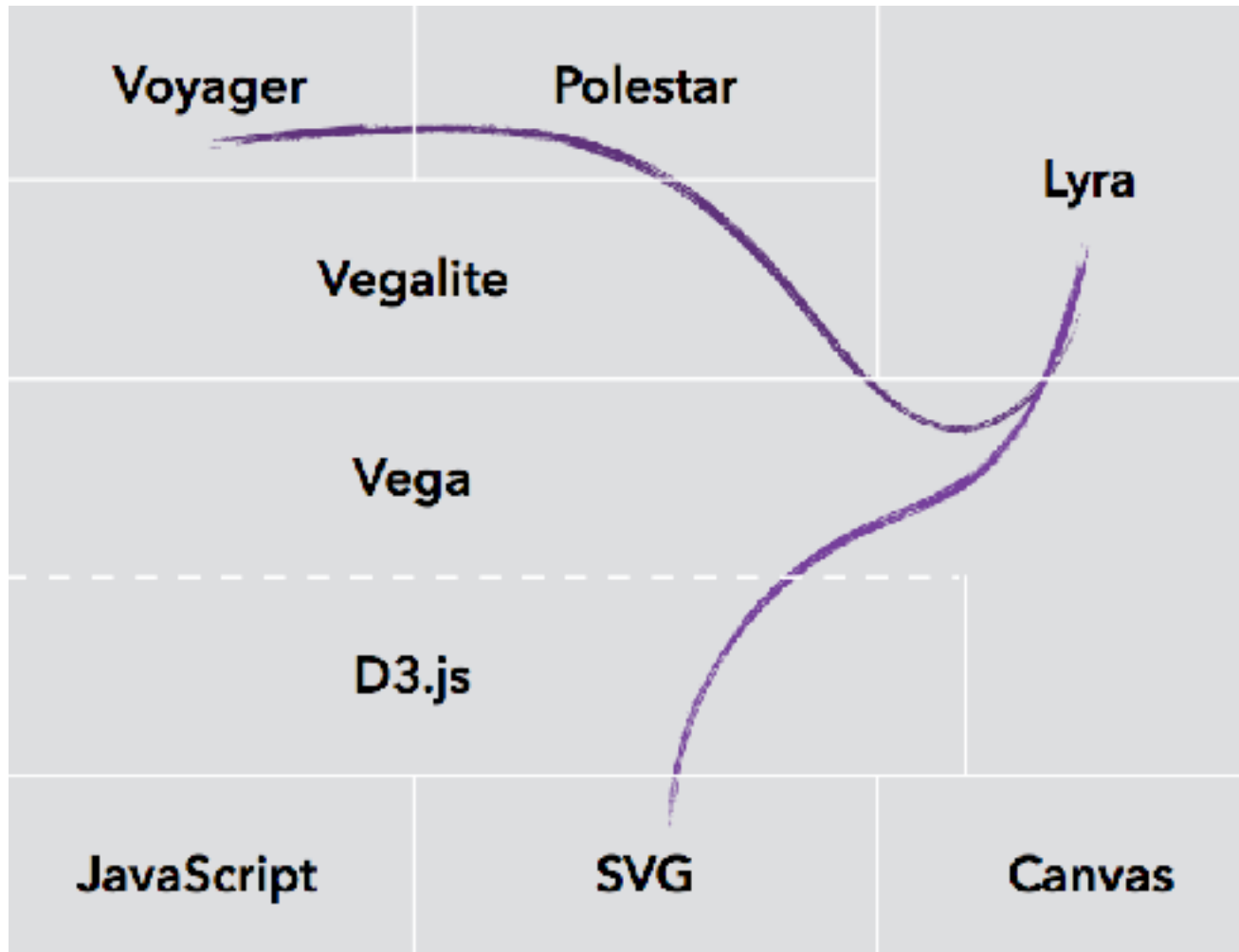
Tools Tradeoffs

- InfoVis-focused
 - Many fundamental techniques built-in
 - Can be faster to get something going
 - Often more difficult to implement something “different”
 - Documentation?
- Generic Graphics
 - More flexible
 - Can customize better
 - Big learning curve
 - Doc is often better
 - Can take a long time to (re)implement basic techniques

Visualization Tools



Many Tools Developed by Prof. Jeffrey Heer, University of Washington



This is just a reference point.
You should try those information
visualization tools out
(optional for those who don't take G53IVP)!

(Optional) Resources

- D3 tutorial: <https://uwdata.github.io/d3-tutorials/>
- Vega tutorial: <https://github.com/vega/vega/wiki/Tutorial>
- Please start working on the course work using R.

Next Lecture

- Topic:
 - Visual Perception

